


- Full 2020
- Home
- Announcements
- Zoom Class Sessions
- Syllabus
- Modules
- Grades
- People
- Course Info
- Library Reserves
- Research Guides
- Canvas Support
- Textbooks
- Mailtool
- NameCoach
- Discussions
- SPS Study Spaces
- Remote Learning Support

## Course Syllabus

[Jump to Today](#)

COLUMBIA UNIVERSITY

School of Professional Studies

### APAN P55430: Applied Text and Natural Language Analytics

Schedule	Dates: September 14, 2020 – December 14, 2020 Time: Mondays, 6:10 PM-8:00 PM Location: Click on "Course Info" for details
Credits	3
Contact Information	<b>Instructor:</b> Javid Huseynov <b>Email</b> jbh2172@columbia.edu <b>Phone</b> 646-726-9526 <b>Office Hours</b> By appointment <b>Response Policy</b> During the term, the easiest way to reach me is via email. You will usually get a response within 48 hours. If you have a question about an assignment, you are advised to email me several days before it is due; if your email arrives within 24 hours of the due date, you may not receive a timely response. <b>Associate: Alon Feshan</b> <b>Email</b> af2981@columbia.edu <b>Zoom</b> Meeting ID: 917 560 0394 Passcode: 299671 <b>Office Hours</b> Wednesdays, 6pm - 7pm EST <b>Response Policy</b> During the term, the easiest way to reach me is via email. You will usually get a response within 48 hours. If you have a question about an assignment, you are advised to email me several days before it is due; if your email arrives within 24 hours of the due date, you may not receive a timely response. <b>Associate: Shabnam Tafreshi</b> <b>Email</b> st3319@columbia.edu <b>Zoom</b> Meeting ID: 220 485 9674 <b>Office Hours</b> Mondays, 4:45pm - 5:45pm EST <b>Response Policy</b> During the term, the easiest way to reach me is via email. You will usually get a response within 48 hours. If you have a question about an assignment, you are advised to email me several days before it is due; if your email arrives within 24 hours of the due date, you may not receive a timely response.
Syllabus Content	
<a href="#">Course Overview</a>   <a href="#">Learning Objectives</a>   <a href="#">Readings</a>   <a href="#">Resources</a>   <a href="#">Course Requirements</a>   <a href="#">Evaluation/Grading</a>   <a href="#">Course Policies</a>   <a href="#">School Policies</a>   <a href="#">Course Schedule</a>	

### Course Overview

With the growth of the Internet in recent decades, there has been an exponential increase of unstructured textual data available from news and social media. This data is invaluable for extracting actionable insights that enhance the scale and the quality of business analytics. But the enormous volume of domain text corpora makes the extraction of meaningful information possible only through the use of advanced natural language processing (NLP) and machine learning techniques. Jobs in the data analysis field increasingly require the use of extracting and analyzing information from diverse sources, structured as well as unstructured. This course will therefore train students in a technology that is seen as an essential part of a data analyst's toolkit.

This course will focus on advanced methods and systems that enable named entity recognition and disambiguation, topic modeling, sentiment analysis, word vector embeddings, abstractive summarization, meaning extraction, and deep learning for NLP. Weekly course lectures will offer a blend of theoretical material and hands-on class exercises, which will be put into practice through weekly assignments. Students who complete the course will be able to practice the gained knowledge as applied NLP-data scientists in various business domains, including sales and marketing, financial modeling, credit risk analysis, legal trust and compliance, intellectual property and contracts management. Some examples include extraction of payment clauses from contracts, establishing customer needs from news, or forecasting industry trends from public announcements.

During the last four weeks of the semester, students will focus on the Term Project that may leverage both unstructured and structured data to discover knowledge about an entity of interest (e.g. company, person, geolocation, etc.). For example, the Term Project may focus on extracting evolving topics or key business events from news articles about a chosen public company, in order to explain changes in its stock price.

The desired outcome of the course is the ability to put the obtained knowledge into practical use. Whether you are taking this course for future academic research, for work in industry, or for an innovative startup idea, this course should help you to master the fundamentals of unstructured data analytics.

**Note:** While extensive programming skills are not required, students are expected to learn, understand and make use of the basic data structures and functions in Python. A refresher session on Python 3 can be offered during one of the weekly classes if necessary. Students have an option to code weekly assignments in either Jupyter Notebook or vanilla Python.

#### Prerequisites

- APAN P55200: Applied Analytics Frameworks and Methods 1
- APAN P55205: Applied Analytics Frameworks and Methods 2

[Back to top](#)

### Learning Objectives

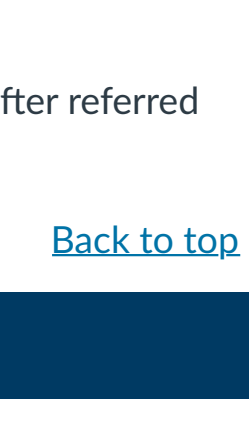
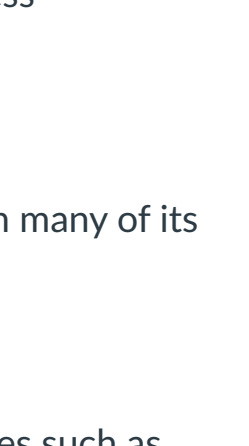
Upon successful completion of this course, you will be able to:

- L1: Derive structured data representation from an unstructured text
- L2: Extract entities and their contexts (relations, keywords, concepts, etc.) from unstructured text
- L3: Model meaning or natural language understanding from text corpora
- L4: Develop descriptive, predictive, and prescriptive analytics models based on text
- L5: Apply topic and sentiment analyses for contextual modeling

[Back to top](#)

### Required Readings

To purchase:

Bengfort, B. (2018). <i>Applied text analysis with Python: Enabling language-aware data products with machine learning</i> e. O'Reilly Media, Inc. ISBN: 1491963042 (Hereafter referred to as "ATAP")	
Hapke, H., Howard, C., & Lane, H. (2019). <i>Natural language processing in action</i> e. Manning Publications. ISBN: 9781617294631 (Hereafter referred to as "NLPA")	

Available online:

- Bird, S., Klein, E., & Loper, E. (2009). *Natural language processing with Python* e. O'Reilly Media, Inc. ISBN: 9780596599707. Also published as *Analyzing text with the natural language toolkit* e. (Hereafter referred to as "NLTK")
- Joshi, P. (2018). *An introduction to text summarization using the TextRank algorithm with Python implementation* e. Analytics Vidhya. Nov. 1, 2018.
- Jurafsky, D., & Martin, J. (2018). *Speech and language processing* e. Third Edition Draft. (Hereafter referred to as "SLP")

[Back to top](#)

### Resources

Columbia University Information Technology

[Columbia University Information Technology](#) e. (CUIT) provides Columbia University students, faculty and staff with central computing and communications services. Students, faculty and staff may access [University-provided and discounted software downloads](#) e.

Columbia University Library

[Columbia's extensive library system](#) e. ranks in the top five academic libraries in the nation, with many of its services and resources available online.

SPS Academic Resources

[The Office of Student Affairs](#) e. provides students with academic counseling and support services such as online tutoring and career coaching.

[Back to top](#)

### Course Requirements (Assignments)

Detailed descriptions and assessment criteria will be provided for all assignments in the Canvas course site.

<b>Term Project (Aligns with L1-L5)</b>  The term project builds on the content and sequence of weekly lectures and assignments. Its purpose is to develop an end-to-end data science solution to a business intelligence problem. The project can be developed by a group or an individual student. Each group can develop own project definition, subject to the instructor's approval.  The proposed data science solution must leverage at least three NLP and text analytics methods studied in the course, e.g. entity recognition, topic modeling, semantic analysis using word embeddings, locality sensitive hashing, etc. Sample problem definitions may include analyzing earnings statements of a company to predict stock price movements, extracting or categorizing contractual clauses, tracking intellectual property infringements from the product announcements, credit and risk assessment as well as sales and market intelligence based on news media.  Groups are expected to implement their proposed end-to-end solution in Python and to present their results in a 10-15 minute live PowerPoint presentation during the final class session. (See "Final Presentation" below for more details).  The Term Project comprises the following deliverables: <ul style="list-style-type: none"><li>• Term Project Proposal</li></ul> In Week 10, you will submit a short (ca. 1 page) proposal in which you select a business use case for the Term Project and define the business impact and value of the project. <ul style="list-style-type: none"><li>• Term Project Deliverable 1: Topic Modeling - Python Implementation<ul style="list-style-type: none"><li>◦ Using topic modeling techniques studied in class, develop a topic taxonomy from the dataset</li><li>◦ Train word vector models based on the dataset and generate topic classifications of articles using semantic similarity</li><li>◦ Aggregate topics over time range to identify their evolution</li></ul></li><li>• Term Project Deliverable 2: Sentiment/Opinion Mining - Python Implementation<ul style="list-style-type: none"><li>◦ Develop sentiment polarity and subjectivity models for the dataset articles using NLTK</li><li>◦ Train a sentiment classifier using Naive Bayes</li></ul></li><li>• Term Project Deliverable 3: Final Project Summary</li></ul> The Final Project Summary shall be presented in the form of either a 3-page written document or 10- to 12-page slide deck and cover the following key elements of the project: <ul style="list-style-type: none"><li>• Background information and objective of the project</li><li>• Data source specification and procurement details</li><li>• Design choices and the rationale for the implemented methodologies</li><li>• Evaluation and model explainability criteria and metrics</li><li>• Findings and conclusions</li></ul>	
<b>Assignments (Aligns with L1-L4)</b>  Weekly assignments give students opportunities to practice foundational skills and are designed to prepare students for a successful completion of the Term Project. <ul style="list-style-type: none"><li>• <b>Assignment 1 (Aligns with L1)</b> : Install PyCharm &amp; Jupyter Notebook and submit the screenshots of installations. Implement an algorithm using Python dictionary, list, and set.</li><li>• <b>Assignment 2 (Aligns with L1)</b>: Write a Python program that:<ul style="list-style-type: none"><li>◦ Implements Webhose.io API calls to obtain a JSON document dataset of web crawls about a chosen entity</li><li>◦ Adds the dataset documents into a persistent storage for use in weekly assignments and the term project</li></ul></li><li>• <b>Assignment 3 (Aligns with L1)</b>: Write a Python program using the regular expression (re), SpaCy and/or NLTK packages to:<ul style="list-style-type: none"><li>◦ Tokenize words and sentences</li><li>◦ Stem and lemmatize work token</li><li>◦ Remove stop words</li><li>◦ List and count n-grams for a given n</li></ul></li><li>• <b>Assignment 4 (Aligns with L2)</b>: Write a Python program to:<ul style="list-style-type: none"><li>◦ Train a simple named entity recognizer using spaCy library</li><li>◦ Recognize and link entities mentioned in text to a knowledge base</li></ul></li><li>• <b>Assignment 5 (Aligns with L2)</b>: Write a Python program to:<ul style="list-style-type: none"><li>◦ Apply TextRank for ranking and selecting key phrases in a document</li><li>◦ Apply LextRank for extractive sentence summarization</li></ul></li><li>• <b>Assignment 6 (Aligns with L3)</b>: Write a Python or PySpark program to train a word or character embedding model using Word2Vec, Glove or FastText algorithms.</li><li>• <b>Assignment 7 (Aligns with L3)</b>: Write a Python program that filters out exactly and/or semantically duplicate articles from the dataset, using SimHash and Word2Vec.</li><li>• <b>Assignment 8 (Aligns with L4)</b>: Write a Python program to calculate document similarity using LSA: matrix decomposition and dimensionality rank reduction.</li><li>• <b>Assignment 9 (Aligns with L5)</b>: Write a Python program that:<ul style="list-style-type: none"><li>◦ Implements LDA topic modeling on the dataset using a choice of Gensim, Scikit-learn or Spark-MLLib libraries</li><li>◦ Extracts topic cluster keywords</li><li>◦ Assigns each article to the relevant cluster based on topics</li></ul></li></ul>	
<b>Final Presentation (Aligns with L1-L5)</b>  At the end of the course, you (or your group) will prepare and deliver a 10-minute presentation on your Term Project. In your presentation, you should address: <ul style="list-style-type: none"><li>• Data source specification &amp; the choice of entity</li><li>• Design choices and implemented solution(s)</li><li>• Evaluation criteria and results</li><li>• Findings and conclusions</li></ul> You will be assessed on how clearly you present your Term Project; how well you facilitate class discussion and Q&A; and on the quality of your slide deck (i.e. does it present information clearly, effectively, and appealingly?). The instructor's frame of reference will be as if you were presenting this Project as an innovative idea to decision-makers within an organization - would the presentation inspire confidence that the activities in your Project are worth undertaking and investing organizational time and resources in?	
<b>Attendance and Participation</b>  Your attendance and active participation are essential to succeed in this course. You are expected to attend all class sessions and come to class prepared, having completed all assigned readings and assignments. During our class meetings, you are expected to engage with your classmates and instructor by answering questions, stating and defending your point of view, and challenging the points of view of others. If you need to miss a class for any reason, please discuss the absence with me in advance.  <i>Final Reflection</i>  Part of your participation grade will be earned through a short (ca. 500 word) reflection that you will write at the end of the course. This reflection will give you the opportunity to address the evolution of your own thinking and learning in this course. You will be asked to identify areas where you were challenged and where you still have lingering questions to explore in the future. You will also be asked to describe how you hope to apply your learnings to your professional life.	

[Back to top](#)

### Evaluation / Grading

The final grade will be calculated as described below:

GRADE CALCULATION		FINAL GRADING SCALE	
ASSIGNMENT	WEIGHT	GRADE	PERCENTAGE
Team Project, including:	50%	A+	98-100 %
• Proposal	10 %	A-	93-97.9 %
• Python Implementation	20%	B+	90-92.9 %
• Final Presentation	20%	B-	87-89.9 %
Assignments	45%	B	83-86.9 %
Final Reflection	5%	B-	80-82.9 %
Attendance and Participation	0%	C+	77-79.9 %
		C	73-76.9 %
		C-	70-72.9 %
		D	60-69.9 %
		F	59.9 % and below

[Back to top](#)

### Course Policies

#### Participation and Attendance

Your attendance and active participation are essential to succeed in this course. You are expected to attend all class sessions and come to class prepared, having completed all assigned readings and assignments. During our class meetings, you are expected to engage with your classmates and instructor by answering questions, stating and defending your point of view, and challenging the points of view of others. If you need to miss a class for any reason, please discuss the absence with me in advance.

#### Late work

Work that is not submitted on the due date noted in the course syllabus without advance notice and permission from the instructor will be graded down 1/3 of a grade for every day it is late (e.g., from a B+ to a B).

#### Citation & Submission

All written assignments must use APA format, cite sources, and be submitted to the course website (not via email).

[Back to top](#)

### School Policies

#### Copyright Policy

Please note—Due to copyright restrictions, online access to this material is limited to instructors and students currently registered for this course. Please be advised that by clicking the link to the electronic materials in this course, you have read and accept the following:

The copyright law of the United States (Title 17, United States Code) governs the making of photocopies or other reproductions of copyrighted materials. Under certain conditions specified in the law, libraries and archives are authorized to furnish a photocopy or other reproduction. One of these specified conditions is that the photocopy or reproduction is not to be "used for any purpose other than private study, scholarship, or research." If a user makes a request for, or later uses, a photocopy or reproduction for purposes in excess of "fair use," that user may be liable for copyright infringement.

#### Academic Integrity

Columbia University expects its students to act with honesty and propriety at all times and to respect the rights of others. It is fundamental University policy that academic dishonesty in any guise or personal conduct of any sort that disrupts the life of the University or denigrates or endangers members of the University community is unacceptable and will be dealt with severely. It is essential to the academic integrity and vitality of this community that individuals do their own work and properly acknowledge the circumstances, ideas, sources, and assistance upon which that work is based. Academic honesty in class assignments and exams is expected of all students at all times.

SPS holds each member of its community responsible for understanding and abiding by the [SPS Academic Integrity and Community Standards](#) e. You are required to read these standards within the first few days of class. Ignorance of the School's policy concerning academic dishonesty shall not be a defense in any disciplinary proceedings.

#### Accessibility

Columbia is committed to providing equal access to qualified students with documented disabilities. A student's disability status and reasonable accommodations are individually determined based upon disability documentation and related information gathered through the intake process. For more information regarding this service, please visit the [University's Health Services website](#) e.

#### Class Recordings

All or portions of the class may be recorded at the discretion of the instructor to support your learning. At any point, the Instructor has the right to discontinue the recording if it is deemed to be obstructive to the learning process.









If the recording is posted, it is considered confidential and it is not acceptable to share the recording outside the purview of the faculty member and registered class.

[Back to top](#)

Course Schedule			
Week	Topic(s)	Readings	Activities / Assignments
1	Introduction & Course Overview	NLTK, <a href="#">Ch 1: Language Processing &amp; Python</a> e (29) ATAP, Ch 1: Language & Computation (18)	• Resources • Session • Assignment 1
2	Data Sources & Crawling	NLTK, <a href="#">Ch 3: Processing Raw Text</a> e (49) ATAP, Ch 2: Building a Custom Corpus (17)	• Resources • Session • Assignment 2
3	Basic Text Processing	NLPA, Ch 2: Build your vocabulary (word tokenization) (58) ATAP, Ch 3: Corpus Preprocessing & Wrangling (17)	• Resources • Session • Assignment 3
4	Information Extraction 1: Named Entity Recognition & Linking	NLPA, Ch 11: Information Extraction (named entity extraction) (25) NLTK, <a href="#">Ch 7: Extracting information from Text</a> e (29)	• Resources • Session • Assignment 4
5	Information Extraction 2: Key Phrase Extraction & Text Summarization	NLPA, Ch. 3: Math with words (33) <a href="#">Joshi, P. (2018). Introduction to text summarization using the TextRank algorithm.</a> e (20)	• Resources • Session • Assignment 5
6	Vector Space Modeling using Neural Networks	NLPA, Ch 6: Reasoning with word vectors (word2vec) (56) SLP, Section 6.8: Vector Semantics: Word2Vec (5)	• Resources • Session • Assignment 6
7	Locality Sensitive Hashing (LSH) & Text Deduplication	NLPA, Appendix F: Locality Sensitive Hashing (11)	• Resources • Session • Assignment 7
8	Topic Modeling	NLPA, Ch 4: Sections 4.1-4.3 (44) ATAP, Ch 6: Clustering for Text Similarity (27)	• Resources • Session • Assignment 8
9	Evolving Topic Classification	No readings	• Resources • Session • Assignment 9
10	Term Project Proposal Presentations	No readings	• Session • Term Project Proposal
11	Deep Learning in NLP	Towards Data Science Blog Post: How Transformers Work BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding (presentation slides by Jacob Devlin, Google AI)	• Resources • Session
12	Sentiment Analysis	No readings	• Session
13	Final Term Project Presentations	No readings	• Session • Term Project Deliverable 1: Live Class Presentation • Term Project Deliverable 2: Python Implementation • Term Project Deliverable 3: Presentation Slide Deck or Summary • Final Reflection

[Back to top](#)

### Course Summary:

Date	Details	
Sun Sep 13, 2020	 <a href="#">Discussion: Gettine Acquainted</a>	to do: 11:59pm
Mon Sep 14, 2020	 <a href="#">1. Session</a>	due by 6:10pm
Mon Sep 21, 2020	 <a href="#">2. Session</a>	due by 11:59pm
Mon Sep 28, 2020	 <a href="#">3. Session</a>	due by 6:10pm
Wed Sep 30, 2020	 <a href="#">Assignment 2</a>	due by 11:59pm
Mon Oct 5, 2020	 <a href="#">4. Session</a>	due by 6:10pm
Wed Oct 7, 2020	 <a href="#">Assignment 3</a>	due by 11:59pm
Mon Oct 12, 2020	 <a href="#">5. Session</a>	due by 6:10pm
Wed Oct 14, 2020	 <a href="#">Assignment 4</a>	due by 11:59pm
Mon Oct 19, 2020	 <a href="#">6. Session</a>	due by 6:10pm
Wed Oct 21, 2020	<a href="#">Assignment 5</a>	due by 11:59pm
Mon Oct 26, 2020	<a href="#">7. Session</a>	due by 6:10pm
Tue Oct 27, 2020	<a href="#">Assignment 6</a>	due by 11:59pm
Wed Nov 4, 2020	<a href="#">Assignment 7</a>	due by 11:59pm
Mon Nov 9, 2020	<a href="#">8. Session</a>	due by 6:10pm
Mon Nov 16, 2020	<a href="#">9. Session</a>	due by 6:10pm
Tue Nov 17, 2020	<a href="#">Assignment 9</a>	due by 11:59pm
Wed Nov 18, 2020	<a href="#">Assignment 8</a>	due by 11:59pm
Mon Nov 23, 2020	<a href="#">10. Session</a>	due by 6:10pm
Sun Nov 29, 2020	<a href="#">Term Project Proposal</a>	due by 11:59pm
Mon Nov 30, 2020	<a href="#">11. Session</a>	due by 6:10pm
Mon Dec 7, 2020	<a href="#">12. Session</a>	due by 6:10pm
	<a href="#">13. Session</a>	due by 6:10pm
Mon Dec 14, 2020	<a href="#">Term Project Deliverable 1: Live Class Presentation</a>	due by 6:10pm
	<a href="#">Final Reflection</a>	due by 11:59pm
Fri Dec 18, 2020	<a href="#">Term Project Deliverable 3: Presentation Slide Deck or Summary</a>	due by 11:59pm
	<a href="#">Term Project Deliverable 2: Python Implementation</a>	due by 11:59pm