

Subject Code	FH6058
Subject Title	Data Analytics for Applied Language Scientists
Credit Value	3
Level	6
Pre-requisite / Co-requisite/ Exclusion	N/A
Objectives	<p>Data analytic skills are increasingly valued in applied language sciences. These include the analysis of natural language (i.e., natural language processing) as well as language-related data (e.g., experimental, survey, demographic data), and the implementation of relevant techniques with programming languages.</p> <p>This subject complements <i>FH6051 Computational Linguistics</i>, which focuses on natural language processing, to further develop competencies in programming, data visualization, and statistical modeling using the <i>Python</i> language. It adopts a thematic and problem-based approach to meet the following objectives.</p> <ul style="list-style-type: none"> • Train students to import, manage, and visualize data relevant to language professionals in <i>Python</i> • Impart working knowledge of basic statistical and machine-learning techniques for language and social research • Develop critical interpretation skills for insights and decision-making • Raise awareness of ethical challenges in today's digital economy, and their implications for language professionals <p>An interactive pedagogical approach will be adopted with balanced assessment tasks and a variety of learning materials. No background in programming, data analytics, or statistics is assumed.</p>
Intended Learning Outcomes (Note 1)	<p>Upon completion of the subject, students will be able to:</p> <ol style="list-style-type: none"> Independently write <i>Python</i> code for research and other work Use key <i>Python</i> libraries like pandas, matplotlib, and statsmodels to manage, visualize, and analyze datasets Understand and use basic statistical and machine learning techniques Apply critical thinking skills to interpret quantitative results for humanistic inquiry Weigh the pros and cons of data analytics along practical and ethical dimensions

<div>Subject Synopsis/ Indicative Syllabus</div> <div>(Note 2)</div>	<div>Introduction</div> <ul style="list-style-type: none">Contemporary data in linguistic and social contextsThe nature and scope of data analytics <div>Data management, visualization, and communication</div> <ul style="list-style-type: none">Spreadsheets and dataframesVisualizing patterns and relationships <div>Data analytics and machine learning</div> <ul style="list-style-type: none">Techniques and methods<ul style="list-style-type: none">General(ized) linear modelMachine learning modelsMonte Carlo simulationsData analytic objectives<ul style="list-style-type: none">Prediction, classification, clustering <div>Data ethics: a critical perspective</div> <ul style="list-style-type: none">Privacy, discrimination, and social inequalities																																								
<div>Teaching/Learning Methodology</div> <div>(Note 3)</div>	<p>Each weekly session will comprise a 2-hour lecture and a 1-hour seminar. During the lecture, the instructor will impart concepts and facilitate/supervise class activities. The seminar, on the other hand, is student-led with groupwork to practice programming skills, analyze data, and present solutions. Students are expected to bring their personal laptops.</p> <p>Teaching and learning is supported by the open-source programming language <i>Python</i> and other online resources. Students should read the prescribed materials and revise previous lessons before each session.</p> <p>There are two individual take-home assignments and an individual project. The take-home assignments require conceptual understanding of data analytic techniques, writing basic programming code, and interpreting the analyses in context. The project requires students to apply these techniques and communicate insights and solutions in simulated professional settings, such as an academic or business conference.</p>																																								
<div>Assessment Methods in Alignment with Intended Learning Outcomes</div> <div>(Note 4)</div>	<table><tr><th rowspan="2">Specific assessment methods/tasks</th><th rowspan="2">% weighting</th><th colspan="5">Intended subject learning outcomes to be assessed (Please tick as appropriate)</th></tr><tr><th>a</th><th>b</th><th>c</th><th>d</th><th>e</th></tr><tr><td>1. Take-home assignment</td><td>35</td><td>✓</td><td>✓</td><td>✓</td><td></td><td></td></tr><tr><td>2. Take-home assignment</td><td>35</td><td>✓</td><td>✓</td><td></td><td>✓</td><td></td></tr><tr><td>3. Individual project</td><td>30</td><td>✓</td><td>✓</td><td>✓</td><td>✓</td><td>✓</td></tr><tr><td>Total</td><td>100 %</td><td colspan="5"></td></tr></table>	Specific assessment methods/tasks	% weighting	Intended subject learning outcomes to be assessed (Please tick as appropriate)					a	b	c	d	e	1. Take-home assignment	35	✓	✓	✓			2. Take-home assignment	35	✓	✓		✓		3. Individual project	30	✓	✓	✓	✓	✓	Total	100 %					
Specific assessment methods/tasks	% weighting			Intended subject learning outcomes to be assessed (Please tick as appropriate)																																					
		a	b	c	d	e																																			
1. Take-home assignment	35	✓	✓	✓																																					
2. Take-home assignment	35	✓	✓		✓																																				
3. Individual project	30	✓	✓	✓	✓	✓																																			
Total	100 %																																								

	<p>The take-home assignments (35% each) require students to demonstrate the intended learning outcomes by combining data analytic and programming with critical thinking skills. The project (30%) requires the same skills with more emphasis on communication and presentation.</p>	
Student Study Effort Expected	Class contact:	
	▪ Lecture (2 hours x 13 weeks)	26 Hrs.
	▪ Seminar (1 hour x 13 weeks)	13 Hrs.
	Other student study effort:	
	▪ Independent reading (3 hours x 13 weeks)	39 Hrs.
	▪ Independent research (2 hours x 13 weeks)	26 Hrs.
	▪ Assignments (2 hours x 13 weeks)	26 Hrs.
	Total student study effort	130 Hrs.

Reading List and References

Recommended text and reference books

Norton, P. et al. (2005). *Beginning Python*. Indianapolis: Wiley.

Bruce, P., Bruce, A., & Gedeck, P. (2020). *Practical Statistics for Data Scientists: 50+ Essential Concepts Using R and Python*. Sebastopol: O'Reilly.

Hai-Jew, S. (Ed.). (2017). *Data Analytics in Digital Humanities*. Cham: Springer.

O'Neil, C. (2016). *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy*. New York: Crown Books.

Tay, D., & Pan, M. X. (Eds.). (2022). *Data Analytics in Cognitive Linguistics: Methods and Insights*. Berlin: De Gruyter Mouton.

Reference papers

Coupé, C. (2018). Modeling linguistic variables with regression models: Addressing non-gaussian distributions, non-independent observations, and non-linear predictors with random effects and generalized additive models for location, scale, and shape. *Frontiers in Psychology*, 9(APR), 1–21.

Kroese, D. P., Brereton, T., Taimre, T., & Botev, Z. I. (2014). Why the Monte Carlo method is so important today. *Wiley Interdisciplinary Reviews: Computational Statistics*, 6(6), 386–392.

Rodriguez, M. Z., Comin, C. H., Casanova, D., Bruno, O. M., Amancio, D. R., Costa, L. da F., & Rodrigues, F. A. (2019). Clustering algorithms: A comparative approach. *PLoS ONE*, 1, 1–34.

Tay, D. (2020). A computerized text and cluster analysis approach to psychotherapy talk. *Language & Psychoanalysis*, 9(1), 1–22.

Python user guides

- Pandas: <https://pandas.pydata.org/>
- Scikit-learn: <https://scikit-learn.org/stable/>
- Seaborn: <https://seaborn.pydata.org/index.html>
- Statsmodels: www.statsmodels.org/stable/user-guide.html
- Numpy: numpy.org/doc/stable/
- Scipy: <https://www.scipy.org/docs.html>

Other online resources

- datacamp.com
- kaggle.com
- towardsdatascience.com