

# Week 1: Introduction

POP77032 Quantitative Text Analysis for Social Scientists

Tom Paskhalis

# Overview

- Module objectives
- Prerequisites and software
- Materials and books
- Module meetings
- Assessment and collaboration
- Weekly schedule

# Module Objectives

- Introduce the fundamentals of working with text as data;
- Extract and prepare textual data for analysis;
- Apply key computational techniques for textual data;
- Practice these concepts using social science examples.

# Module Materials

- Module website: [tom.paskhal.is/POP77032](http://tom.paskhal.is/POP77032)
- Blackboard

# Books

- Justin Grimmer, Margaret E. Roberts, and Brandon M. Stewart. 2022. *Text as Data: A New Framework for Machine Learning and the Social Sciences*. Princeton, PA: Princeton University Press
- Daniel Jurafsky and James H. Martin. 2026. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition with Language Models*. 3rd ed. Draft.

Also:

- Christopher Manning and Hinrich Schütze. 1999. *Foundations of Statistical Natural Language Processing*. The MIT Press
- Jacob Eisenstein. 2019. *Introduction to Natural Language Processing*. Cambridge, MA: The MIT Press.
- Klaus Krippendorff. 2019. *Content Analysis: An Introduction to Its Methodology*. 4th ed. Thousand Oaks, CA: SAGE Publications

# Additional Online Materials

- quanteda
- APIs for Social Scientists: A Collaborative Review
- Text Mining with R

# Prerequisites and Software

- Intermediate module - familiarity with basic statistical concepts and programming in R/Python is assumed.
- Laptop with Windows/Mac/Linux OS (no Chrome books)
- Required software:
  - **Jupyter** - web-based interactive computational environment
  - **Python** (version 3+) - versatile programming language
  - **R** (version 4+) - statistical programming language
- Additional software:
  - **JupyterLab Desktop** - desktop application for Jupyter Notebooks
  - **RStudio** - integrated development environment for R
  - **Spyder** - integrated development environment for Python
  - **Visual Studio Code** - feature-rich text editor

# Module Meetings

- 2-hour lecture
  - Until Reading Week - Wednesday 16:00-18:00 in 4050A Arts Building
  - After Reading Week - Wednesday 16:00-18:00 in 5052 Arts Building
- 2-hour tutorials
  - Until Reading Week - Friday 14:00-16:00 in AP0.09 Aras an Phiarsaigh
  - After Reading Week - Friday 14:00-16:00 in 1.24 D'Olier Street
- Office hours:
  - Friday 11:00 - 13:00 online or in-person (booking required)

# Assessment

- 3 programming exercises (40%)
- Research paper (60%)
  - Approximately 10 pages and 3,000-4,000 words (references excluded)
  - Due by **23:59 Wednesday, 22 April 2026**

# Plagiarism Policy

- Plagiarising computer code is as serious as plagiarising text (see Google LLC v. Oracle America, Inc.)
- All submitted programming assignments and final project should be done individually;
- You may discuss general approaches to solutions with your peers;
- But do not share or view each others code;
- You can use online resources but give credit in the comments.

# Generative AI Policy

- The use of generative AI is permitted.
- However:
  - No part of the module content can be used in a prompt;
  - It needs to be explicitly acknowledged in the submission;
  - You need to state the version of the model used.
- Hardware permitting, I recommend using local offline models installed on your machine.
- E.g. check LM Studio as a user-friendly interface to different models.

# Module Outline

<b>Week</b>	<b>Date</b>	<b>Topic</b>	<b>Released</b>	<b>Due</b>
1	21 January	Introduction		
2	28 January	Words and Tokens		Assignment 1
3	4 February	Quantifying Texts		
4	11 February	Dictionaries and Sentiment		Assignment 1
5	18 February	Supervised Modelling		Assignment 2
6	25 February	Unsupervised Modelling		
7	4 March	-		Assignment 2
8	11 March	Beyond Bag-of-Words		
9	18 March	Embeddings		Assignment 3
10	25 March	Neural Networks		
11	1 April	Transformers		Assignment 3
12	8 April	Large Language Models		

# Next

- Introduction to QTA