

課程概述與目標

大數據 (Big Data) 的時代, 數位資料累積與增長的速度已經遠遠倍增於人類史上的任何階段。這樣一種鉅量資料風潮, 不僅改變了人文社會與自然科學研究的面貌, 在各項產業也產生了分析資料輔佐決策的迫切需求。在此背景下, 數位素養 (digital literacy) 已經成為現代公民必須具備的基本素養之一, 近年來新興的**資料科學家** (data scientist) 更成為當前最為熱門的行業之一。然而由於數據資料的發展, 隨著社交媒體與社會網路的發展, 非結構性的文本資料所佔比例已經遠超過結構性的表格性資料, 使得文本的語言分析在資料科學發展中的角色顯得愈來愈重要, 特別對於人文與社會科學的學生而言, 更是開啓了一個結合數位科技與人文關懷的新的發展方向。本課程的設計, 就是在這個動機之下, 透過介紹**語言學與文本分析** 知識結合統計計算, 並運用新的教學與實作平台, 希望能夠刺激人文、社會、傳播與其他財務管理、醫學等各領域學生之間的互動協作與學習, 培養跨領域的興趣與分析能力。

資料科學家的工作, 可以視為是一個探索、預測與解讀資料意義的互動歷程。而語言分析的工作, 在了解文本資料的語意與情緒表現上是重要的關鍵。本課程結合了目前統計程式設計與自然語言處理技術, 以較為簡潔容易入門的設計與實際操作導引, 希望可以讓毫無相關程式學習基礎的學生在本課程的帶領下, 達到以下的學習目標:

- 了解結構與非結構性資料的特性與預處理工作, 特別是針對中文文本中呈現的語言特性的處理方法。
- 了解中文的語言特性與文本解析 (text analytics) 的基本概念。
- 選擇適當的變數與特徵並加以合理調製, 對之進行描述統計與視覺探勘, 針對不同的問題點與數據類型, 找出適當的圖形表達與統計分析。
- 學習簡易的自然語言處理與機器學習預測模式, 並應用在自己關心的領域。

Syllabus (tentative)

Week	Date	Topic	Lab
1	09/17	Orientation	
2	09/24	Introduction to Data Science and Text Analytics	
3	10/01	Introduction to Data Science and Text Analytics	
4	10/08	Preparing / Preprocessing text and linguistics ABC	
5	10/15	Preparing / Preprocessing text and linguistics ABC	
6	10/22	Exploratory data analysis and Infographics	
7	10/29	Exploratory data analysis and Infographics	
8	11/05	Corpus and natural language processing	
9	11/12	Corpus and natural language processing	
10	11/19	Text classification and clustering	
11	11/26	mini-Hackathon [開始規劃期末展演]	
12	12/03	Text classification and clustering	
13	12/10	Topic modeling	
14	12/17	Sentiment analysis	
15	12/24	Stylometrics and personality detection	
16	12/31	Discussion	
17	01/07	Term project presentation (un-conference)	
18	01/14	Final term project report due	

課程概述與目標

- 上課方式 分爲兩部分：針對每週主題有 2 個小時的課堂講解，與 1 小時的實作教學與練習。每次上課都會分派作業，以階段性的題組一步一步的累積相關知識。每次作業預估的工作量約每週 5 小時。所有學習活動的進度、歷程與成果都在網路平台上進行，一方面老師與助教容易掌握，同學之間也可彼此觀摩學習。
- 成績評量項目與方式：課後作業與課堂表現 (30%)；期中評量 (20%)；期末計畫展演、組織與報告 (50%)。

期中評量將以黑客松 (hackathon) 方式進行公開之跨領域小組協作競賽。期末計畫的部分，則將邀請語言科技新創公司團隊於期中之後出題，進行計畫實作練習。期末計畫展演時則邀請該團隊 CEO 來參與評分，也增進同學對於如何進行知識應用的了解。計畫展演的方式將以開放的短會 unconference 方式進行，由各組同學彼此協調與自行組織，藉此練習學習、組織與團隊協作的翻轉學習精神。

由於此課程內容較新，目前並未有完整的中文教材，此外，考量到本課程是以**推廣通識精神**爲目的，在閱讀教材的設計上，將以授課教師自己編撰的教材講義爲主，輔以列在延伸閱讀中的英語延伸教材。由於本課程著重在「觀念激盪」與「創新實作」。有別於傳統上以紙本書籍閱讀爲主的課程設計，本課程將以預先設計的「主題」爲核心，透過翻轉式的「網路閱讀」爲學習方式，鼓勵同學在此過程中，學習到協力合作並累積自學經驗，結合自身專業而能轉化成未來可能的應用。所設計的「主題」包括：「文本分類」、「情緒分析」、「立場偵測」、「個人文體風格」、「語言模式抽取與分析」(語意語用與言談等等)。