



THE UNIVERSITY  
of EDINBURGH

## Text Technologies for Data Science

INFR11145

# Comparing Text Corpora

Instructor:

Steve Wilson

Assignment Project Exam Help

11-Nov-2020

1

<https://powcoder.com>

## Initial Text Analysis

Add WeChat powcoder

- Scenario: you are given access to a new dataset
  - 2 corpora, each contains thousands of plain text files
  - You want to understand and quantify:
    - What is the *content* of these documents? What are they *about*?
    - How does the content of these corpora *differ*?
- What are some things you might try first?



2

## Lecture Objectives

- Analyze text corpora
  - Content analysis background
  - Word-level differences
  - Dictionaries and Lexicons
  - Topic modeling
  - Annotation + classification

Assignment Project Exam Help

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

3

<https://powcoder.com>

## Content Analysis

Add WeChat powcoder

- Goal: given some documents determine
  - What are the types of content present? (themes/topics)
  - Which documents contain which topics?
- Traditionally a manual process
  1. Read a subset of documents, define themes/topics
  2. Determine consistent coding\* methodology
  3. Read all documents and label them according to codes
  4. Check agreement between human coders
  5. Settle disagreements via a third-party
  6. Analyze resulting annotations

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

4

## Content Analysis

- Can this process be automated?
  - Yes, to an extent
- *Should* this process be automated?
  - Humans are better than machines at this task (for now?)
  - Computers are *much, much* faster
    - Avg. human reading speed: 250 wpm
    - Assume 1K words/document, 50K documents...
      - Average person needs > 4 months to read
      - This is a **relatively small** corpus for modern NLP
    - Modern computers can process millions of words/second

## Assignment Project Exam Help

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

5

<https://powcoder.com>

## Automated Content Analysis

- |                           |   |                            |
|---------------------------|---|----------------------------|
| • Single corpus/class     |   | • Multiple corpora/classes |
| • Word frequency analysis | ↔ | • Word-level differences   |
| • Dictionaries & Lexicons | ↔ | • Dominance Scores         |
| • Topic modelling         | ↔ | • Topic-level differences  |

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

6

THE UNIVERSITY  
of EDINBURGH

<https://powcoder.com>

THE UNIVERSITY  
of EDINBURGH

## Word-level Differences

- Which words best characterize a corpus?
  - Need a reference corpus
- Some methods to do this:
  - Mutual information
  - Chi squared
- Can also be used for *feature selection*

## Assignment Project Exam Help

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

9

<https://powcoder.com>

## Mutual Information

- $I(X;Y)$ 
  - How much can I learn about X by observing Y?
  - Is the same as *information gain*
  - Is **not** the same as *pointwise mutual information*
- We want to learn about important words in our corpus
- What should X and Y be?
  - $X = U$  = document contains term t (Boolean)
  - $Y = C$  = class is the target class (Boolean)

$$I(U;C) = \sum_{e_t \in \{1,0\}} \sum_{e_c \in \{1,0\}} P(U = e_t, C = e_c) \log_2 \frac{P(U = e_t, C = e_c)}{P(U = e_t)P(C = e_c)}$$

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

10

Add WeChat powcoder

## Mutual Information

$$I(U;C) = \sum_{e_t \in \{1,0\}} \sum_{e_c \in \{1,0\}} P(U = e_t, C = e_c) \log_2 \frac{P(U = e_t, C = e_c)}{P(U = e_t)P(C = e_c)}$$

- Given count data for 2 classes, can be computed as:

$$I(U;C) = \frac{N_{11}}{N} \log_2 \frac{NN_{11}}{N_{1.}N_{.1}} + \frac{N_{01}}{N} \log_2 \frac{NN_{01}}{N_{0.}N_{.1}} \\ + \frac{N_{10}}{N} \log_2 \frac{NN_{10}}{N_{1.}N_{.0}} + \frac{N_{00}}{N} \log_2 \frac{NN_{00}}{N_{0.}N_{.0}}$$

## Assignment Project Exam Help

Source: Manning, Raghavan, and Schütze, 2008

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

11

<https://powcoder.com>

## Mutual Information

$$I(U;C) = \frac{N_{11}}{N} \log_2 \frac{NN_{11}}{N_{1.}N_{.1}} + \frac{N_{01}}{N} \log_2 \frac{NN_{01}}{N_{0.}N_{.1}} \\ + \frac{N_{10}}{N} \log_2 \frac{NN_{10}}{N_{1.}N_{.0}} + \frac{N_{00}}{N} \log_2 \frac{NN_{00}}{N_{0.}N_{.0}}$$

- Example:
  - What is  $I(U;C)$  given these values?

	$e_c = e_{poultry} = 1$	$e_c = e_{poultry} = 0$
$e_t = e_{export} = 1$	$N_{11} = 49$	$N_{10} = 27,652$
$e_t = e_{export} = 0$	$N_{01} = 141$	$N_{00} = 774,106$

Example: Manning, Raghavan, and Schütze, 2008

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

12

## Mutual Information for News Data

UK		China		poultry	
london	0.1925	china	0.0997	poultry	0.0013
uk	0.0755	chinese	0.0523	meat	0.0008
british	0.0596	beijing	0.0444	chicken	0.0006
stg	0.0555	yuan	0.0344	agriculture	0.0005
britain	0.0469	shanghai	0.0292	avian	0.0004
plc	0.0357	hong	0.0198	broiler	0.0003
england	0.0238	kong	0.0195	veterinary	0.0003
pence	0.0212	xinhua	0.0155	birds	0.0003
pounds	0.0149	province	0.0117	inspection	0.0003
english	0.0126	taiwan	0.0108	pathogenic	0.0003

coffee		elections		sports	
coffee	0.0111	election	0.0519	soccer	0.0681
bags	0.0042	elections	0.0342	cup	0.0515
growers	0.0025	polls	0.0339	match	0.0441
kg	0.0019	voters	0.0315	matches	0.0408
colombia	0.0018	party	0.0303	played	0.0388
brazil	0.0016	vote	0.0299	league	0.0386
export	0.0014	poll	0.0225	beat	0.0301
exporters	0.0013	candidate	0.0202	game	0.0299
exports	0.0013	campaign	0.0202	games	0.0284
crop	0.0012	democratic	0.0198	team	0.0264

Example: Manning, Raghavan, and Schütze, 2008

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

Assignment Project Exam Help

14

<https://powcoder.com>

## Chi-squared

Add WeChat powcoder

- Hypothesis testing approach
- $H_0$ : Term appearance is independent from a document's class
  - i.e.,  $P(U=1, C=1) = P(U=1)P(C=1)$
- Compute:

$$X^2(\mathbb{D}, t, c) = \sum_{e_t \in \{0,1\}} \sum_{e_c \in \{0,1\}} \frac{(N_{e_t e_c} - E_{e_t e_c})^2}{E_{e_t e_c}}$$

- Or to directly plug in values like before:

$$X^2(\mathbb{D}, t, c) = \frac{(N_{11} + N_{10} + N_{01} + N_{00}) \times (N_{11}N_{00} - N_{10}N_{01})^2}{(N_{11} + N_{01}) \times (N_{11} + N_{10}) \times (N_{10} + N_{00}) \times (N_{01} + N_{00})}$$

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

15

## Chi-squared

$$X^2(\mathbb{D}, t, c) = \frac{(N_{11} + N_{10} + N_{01} + N_{00}) \times (N_{11}N_{00} - N_{10}N_{01})^2}{(N_{11} + N_{01}) \times (N_{11} + N_{10}) \times (N_{10} + N_{00}) \times (N_{01} + N_{00})}$$

- Example

- What is the value of  $X^2$  given the example data?

	$e_c = e_{poultry} = 1$	$e_c = e_{poultry} = 0$
$e_t = e_{export} = 1$	$N_{11} = 49$	$N_{10} = 27,652$
$e_t = e_{export} = 0$	$N_{01} = 141$	$N_{00} = 774,106$

## Assignment Project Exam Help

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

16

<https://powcoder.com>

Add WeChat powcoder

## Dictionaries and Lexicons

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

17



## Dictionaries and Lexicons

- What if we know what we are looking for?
- Dictionaries (lexicons) are prebuilt mappings
  - Category -> word list
  - E.g., a tiny sentiment lexicon:
    - Positive: good, great, happy, amazing, wonderful, best, incredible
    - Negative: terrible, horrible, bad, awful, nasty, gross, worst, poor
- Domain can be important
  - “**unpredictable** movie plot” ✓
  - “**unpredictable** coffee pot” ✗

## Assignment Project Exam Help

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

18

<https://powcoder.com>

## Dictionaries and Lexicons

- How to get a score per category?

$$\frac{\text{num\_dictionary\_words\_in\_document}}{\text{num\_total\_words\_in\_document}}$$

- That's it!
- Can also be used as machine learning features
- A more advanced approaches to quantifying categories (optional reading)
  - <https://www.ncbi.nlm.nih.gov/pubmed/28364281>

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

19

## Some Dictionaries

- LIWC (Pennebaker et al. 2015)
- General Inquirer (Stone 1997)
- Roget's Thesaurus Categories
- VADER (Hutto and Gilbert, 2014)
- Sentiwordnet (Esuli and Sebastiani 2006)
- Wordnet Domains (Magnini and Cavaglia, 2000)
- EmoLex (Mohammad and Turney, 2010)
- Empath (Fast et al., 2016)
- Personal Values Lexicon (Wilson et al., 2018)
- ...

## Assignment Project Exam Help

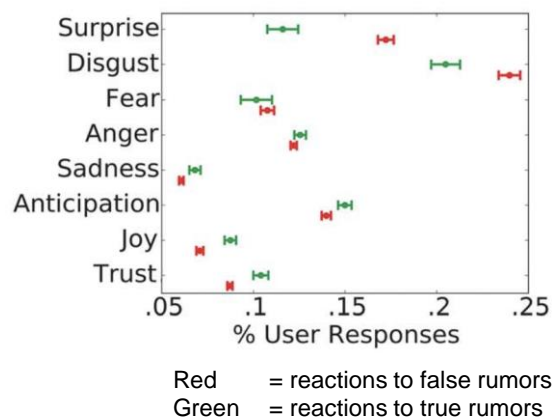
Steve Wilson, TTDS 2020/2021



20

<https://powcoder.com>

## Reactions to Rumor Tweets With EmoLex



Vosoughi, Roy, and Aral, 2018

Steve Wilson, TTDS 2020/2021



21

## Dominance Scores

- The dominance score for a category w.r.t. a corpus:

$$\frac{\text{category\_score\_in\_target\_corpus}}{\text{category\_score\_in\_background\_corpus}}$$

Assignment Project Exam Help

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

22

<https://powcoder.com>

## LIWC category dominance scores

Truthful				Deceptive			
Interviews		Trials		Interviews		Trials	
Class	Score	Class	Score	Class	Score	Class	Score
Metaphor	2.98	You	3.99	Assent	4.81	Anger	2.61
Money	2.74	Family	3.07	Past	2.59	Anxiety	2.61
Inhibition	2.74	Home	2.45	Sexual	2.00	Certain	2.28
Home	2.13	Humans	1.87	Other	1.87	Death	1.96
Humans	2.02	Posemo	1.81	Motion	1.68	Physical	1.77
Family	1.96	Insight	1.64	Negemo	1.44	Negemo	1.52

Pérez-Rosas et al, 2015

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

23

## Topic Level Analysis

# Assignment Project Exam Help

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

24

<https://powcoder.com>

## Intro to Topic Modelling

- Goals are similar to traditional content analysis:
  - What are the main themes/topics in this corpus?
  - Which documents contain which topics?

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

25

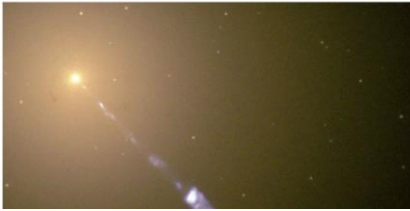
Add WeChat powcoder

## Topic Models

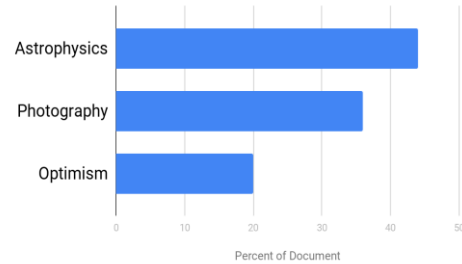
The New York Times

### ***Expected Soon: First-Ever Photo of a Black Hole***

Have astronomers finally recorded an image of a black hole? The world will know on Wednesday.



Topic Distribution



# Assignment Project Exam Help 26

Steve Wilson, TTDS 2020/2021

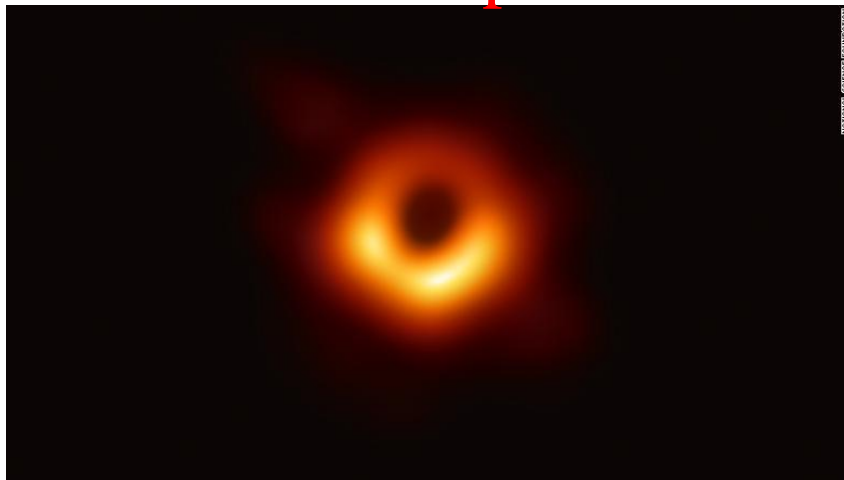


THE UNIVERSITY  
of EDINBURGH

26

<https://powcoder.com>

## Add WeChat powcoder



27

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

27

human	evolution	disease	computer
genome	evolutionary	host	models
dna	species	bacteria	information
genetic	organisms	diseases	data
genes	life	resistance	computers
sequence	origin	bacterial	system
gene	biology	new	network
molecular	groups	strains	systems
sequencing	phylogenetic	control	model
map	living	infectious	parallel
information	diversity	malaria	methods
genetics	group	parasite	networks
mapping	new	parasites	software
project	two	united	new
sequences	common	tuberculosis	simulation

Example from David Blei

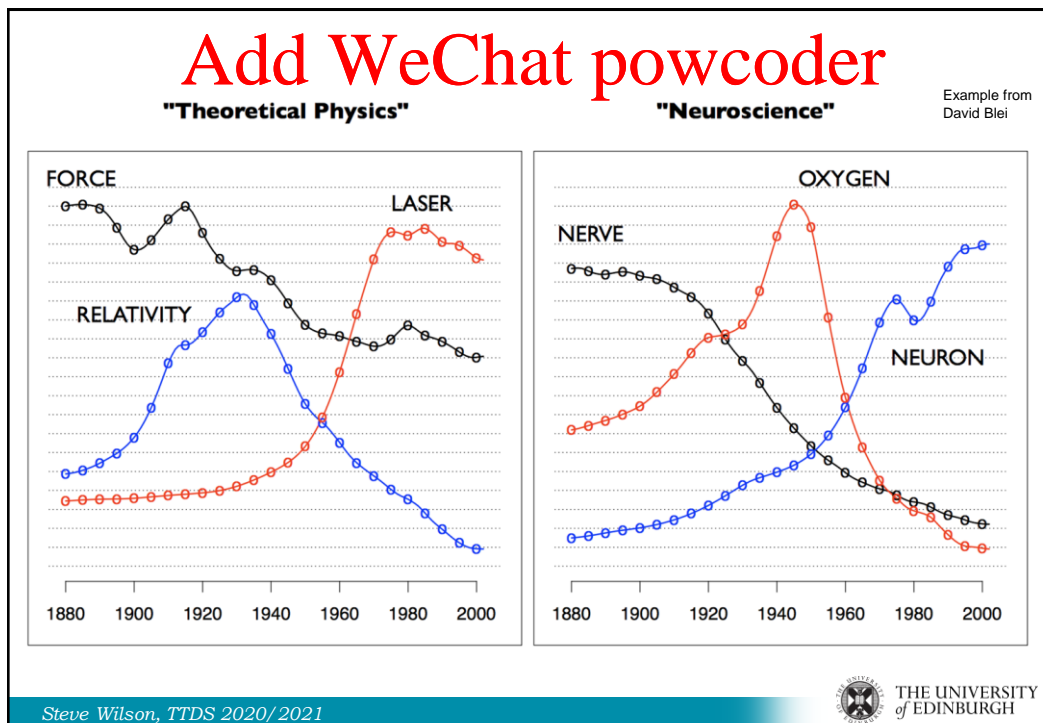
28

Steve Wilson, TTDS 2020/2021

THE UNIVERSITY of EDINBURGH

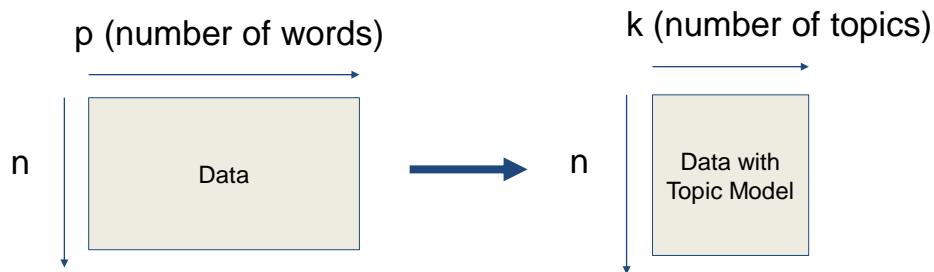
28

<https://powcoder.com>



29

## Dimensionality Reduction



Assignment Project Exam Help 30

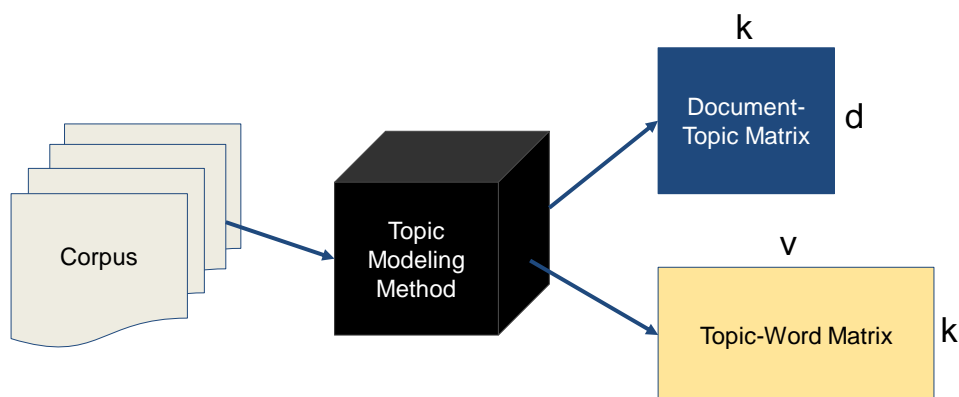
Steve Wilson, TTDS 2020/2021



30

<https://powcoder.com>

## Topic Modeling



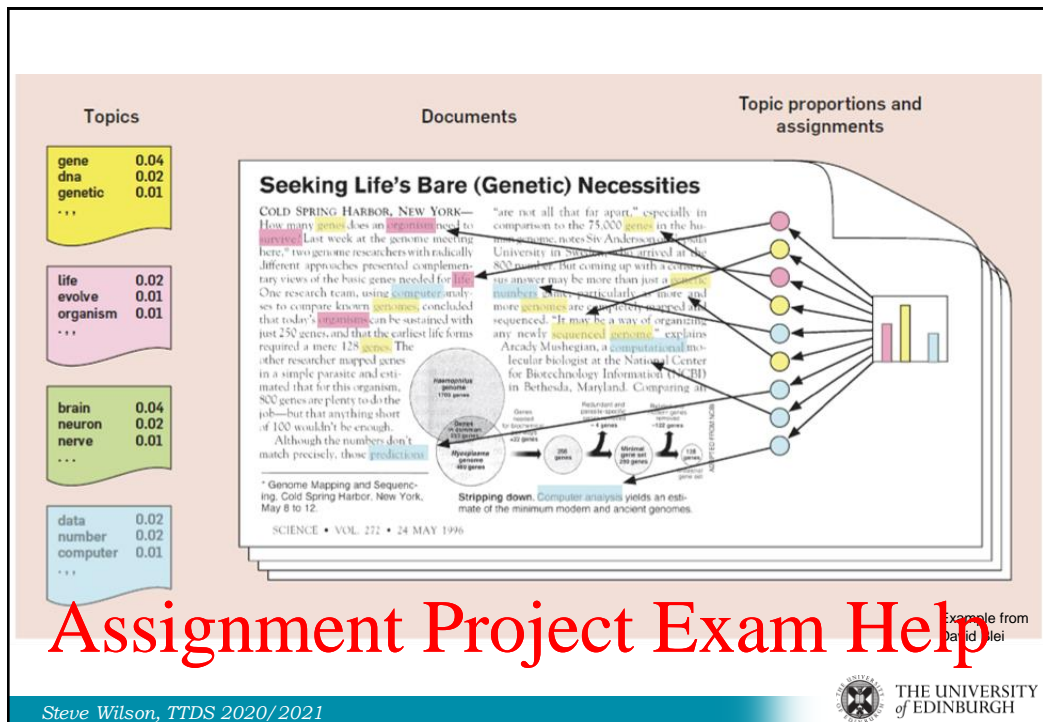
31

Steve Wilson, TTDS 2020/2021



31

Add WeChat powcoder



32

<https://powcoder.com>

## Topic Models Add WeChat powcoder

- Most often used for text data, but can also be applied in other settings:
  - Bioinformatics (Liu et al. 2016)
  - Computer code (McBurney et al. 2014)
  - Music (Hu and Saul 2009)
  - Network data (Cha and Cho 2014)

33

Steve Wilson, TTDS 2020/2021

33



## Topic Modeling Methods

- Most popular: Latent Dirichlet Allocation (LDA)
  - Introduced by David Blei, Andrew Ng, and Michael Jordan (2003)
- Other methods include
  - pLSI
  - PCA-based methods
  - Non-negative matrix factorization
  - Deep learning based topic modeling
  - ...

Assignment Project Exam Help 34

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

34

<https://powcoder.com>

## Topic Modeling Methods

- Most popular: Latent Dirichlet Allocation (LDA)
  - Introduced by David Blei, Andrew Ng, and Michael Jordan (2003)
- Other methods include
  - pLSI
  - PCA-based methods
  - Non-negative matrix factorization
  - Deep learning based topic modeling
  - ...

35

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

35

## Latent Dirichlet Allocation (LDA)

- More details coming up in next lecture...

Assignment Project Exam Help 36

Steve Wilson, TTDS 2020/2021



THE UNIVERSITY  
of EDINBURGH

36

<https://powcoder.com>

Add WeChat powcoder