



THE UNIVERSITY
of EDINBURGH

Text Technologies for Data Science

INFR11145

Comparing Text Corpora I

Instructor:
Steve Wilson

13-Nov-2019

1

Pre-Lecture

- Today
 - Lecture: Comparing Text Corpora

Steve Wilson, TTDS 2019/2020



2

1

Initial Text Analysis

- 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?

Steve Wilson, TTDS 2019/2020



3

Lecture Objectives

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

Steve Wilson, TTDS 2019/2020



5

2

Content Analysis

- 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 2019/2020



6

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

Steve Wilson, TTDS 2019/2020



7

3

Automated Content Analysis

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

Steve Wilson, TTDS 2019/2020



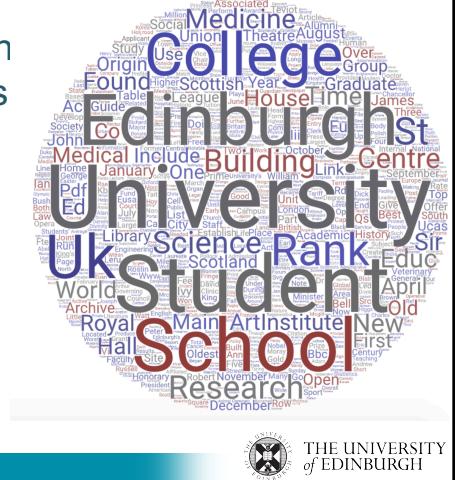
Word Level Analysis

Steve Wilson, TTDS 2019/2020



Word frequency analysis

- Very simple starting point
 1. Preprocess as usual (lowercasing? stemming?...)
 2. Count words
 3. Normalize by document length
 4. Average across all documents



Steve Wilson, TTDS 2019/2020



THE UNIVERSITY
of EDINBURGH

10

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*

Steve Wilson, TTDS 2019/2020



THE UNIVERSITY
of EDINBURGH

11

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 2019/2020



12

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, can be computed as:

$$\begin{aligned} I(U;C) &= \frac{N_{11}}{N} \log_2 \frac{NN_{11}}{N_{1\cdot}N_{\cdot1}} + \frac{N_{01}}{N} \log_2 \frac{NN_{01}}{N_{0\cdot}N_{\cdot1}} \\ &\quad + \frac{N_{10}}{N} \log_2 \frac{NN_{10}}{N_{1\cdot}N_{\cdot0}} + \frac{N_{00}}{N} \log_2 \frac{NN_{00}}{N_{0\cdot}N_{\cdot0}} \end{aligned}$$

Source: Manning, Raghavan, and Schütze, 2008

Steve Wilson, TTDS 2019/2020



13

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?

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

Example: Manning, Raghavan, and Schütze, 2008

Steve Wilson, TTDS 2019/2020



14

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 2019/2020



16

Chi-squared

- 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 2019/2020



17

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?

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

Steve Wilson, TTDS 2019/2020



18

Dictionaries and Lexicons

Steve Wilson, TTDS 2019/2020



20

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” ✗

Steve Wilson, TTDS 2019/2020



21

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 2019/2020



22

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)
- ...

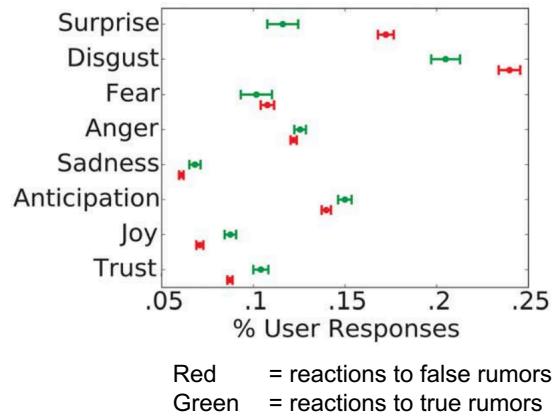
Steve Wilson, TTDS 2019/2020



23

10

Reactions to Rumor Tweets with EmoLex



Vosoughi, Roy, and Aral, 2018

Steve Wilson, TTDS 2019/2020



24

Dominance Scores

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

$$\frac{\text{category_score_in_target_corpus}}{\text{category_score_in_background_corpus}}$$

- From Mihalcea and Pulman, 2009

Steve Wilson, TTDS 2019/2020



25

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 2019/2020



26

Topic Level Analysis

Steve Wilson, TTDS 2019/2020



27

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 2019/2020



28

Topic Models

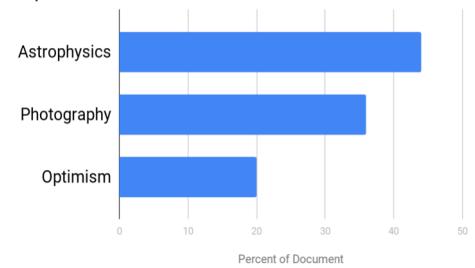
The New York Times

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



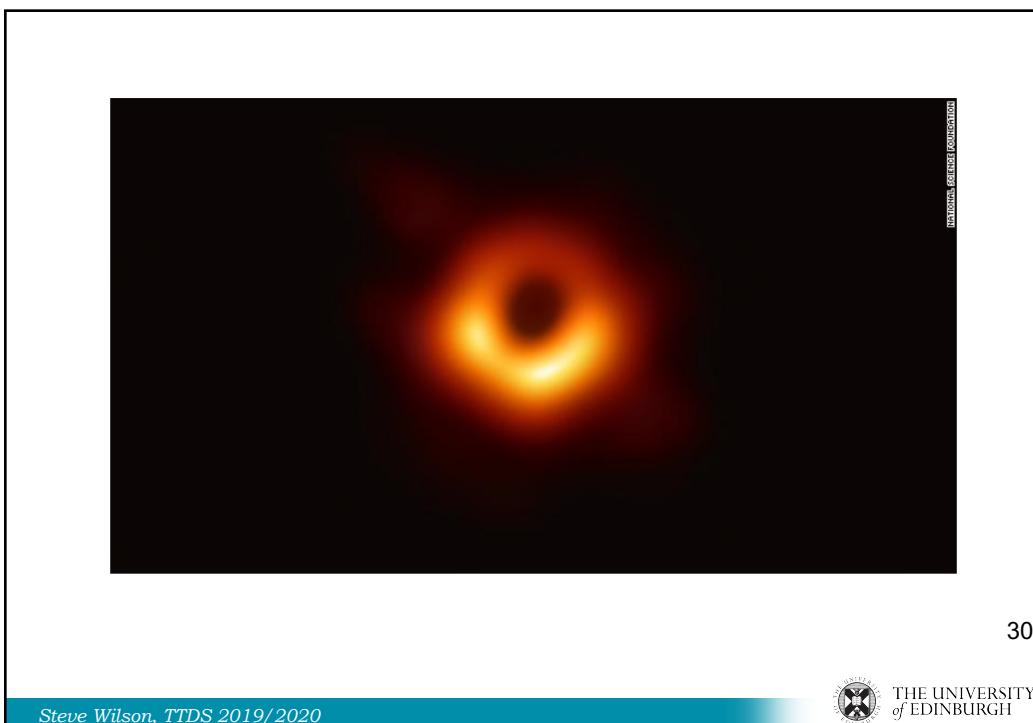
Steve Wilson, TTDS 2019/2020



29

29

13



30

Steve Wilson, TTDS 2019/2020



30

Topic Models

- 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)

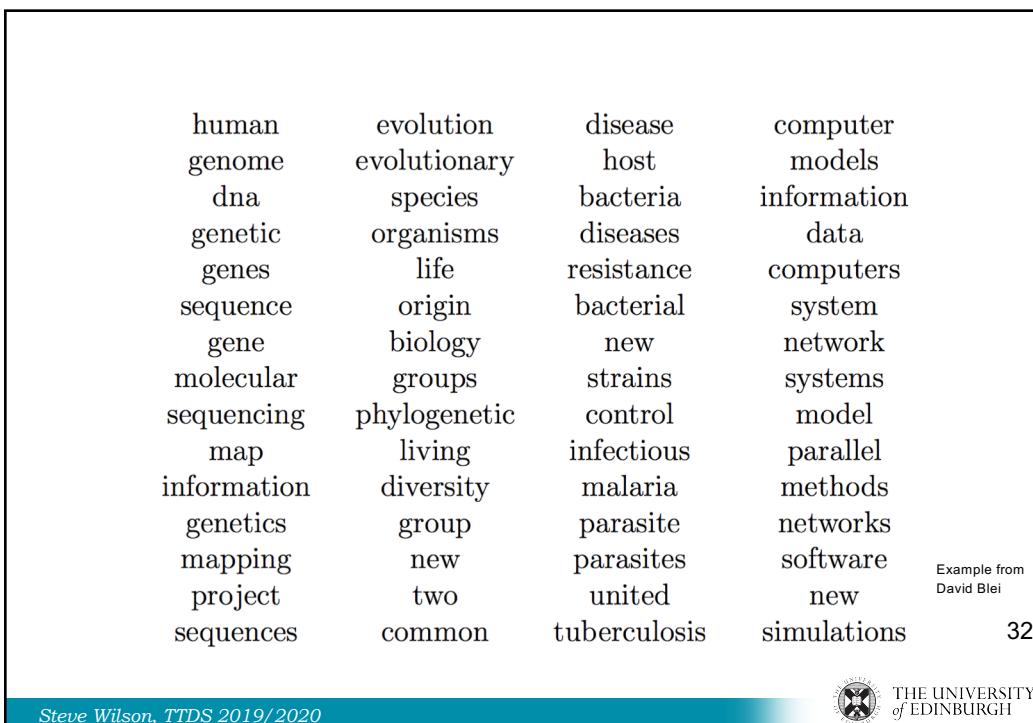
31

Steve Wilson, TTDS 2019/2020



31

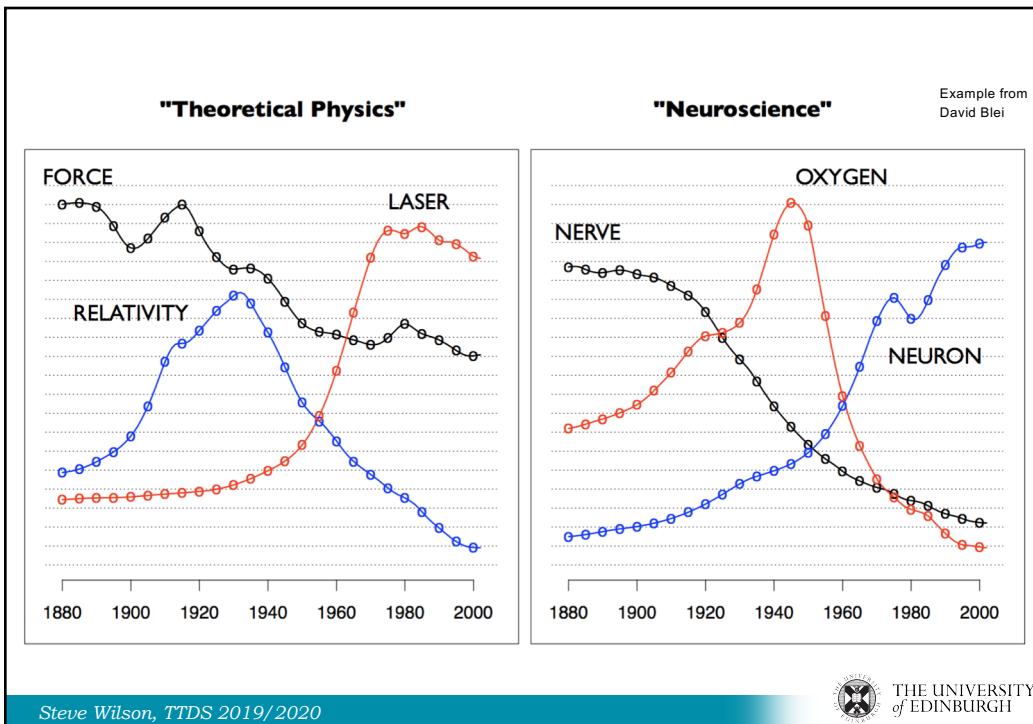
14



Steve Wilson, TTDS 2019/2020

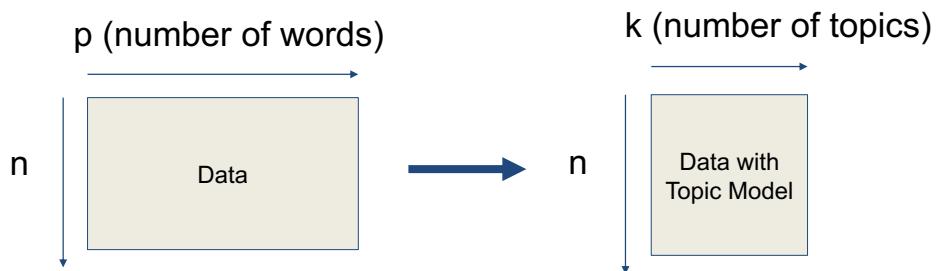


32



33

Dimensionality Reduction



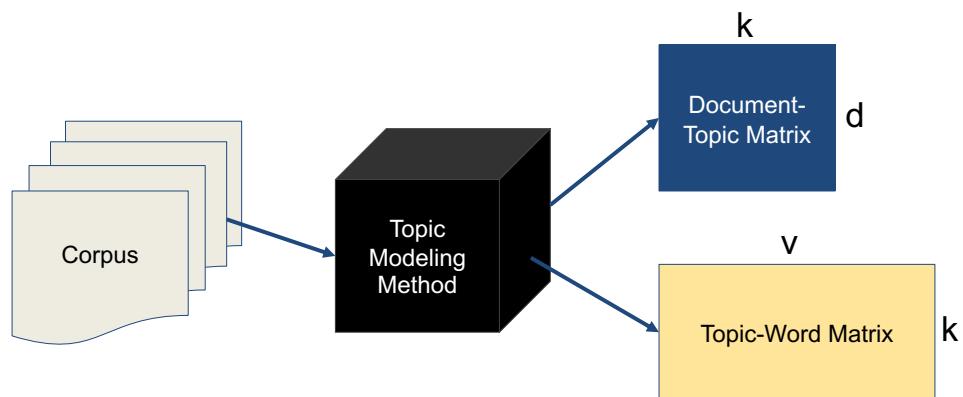
34

Steve Wilson, TTDS 2019/2020



34

Topic Modeling



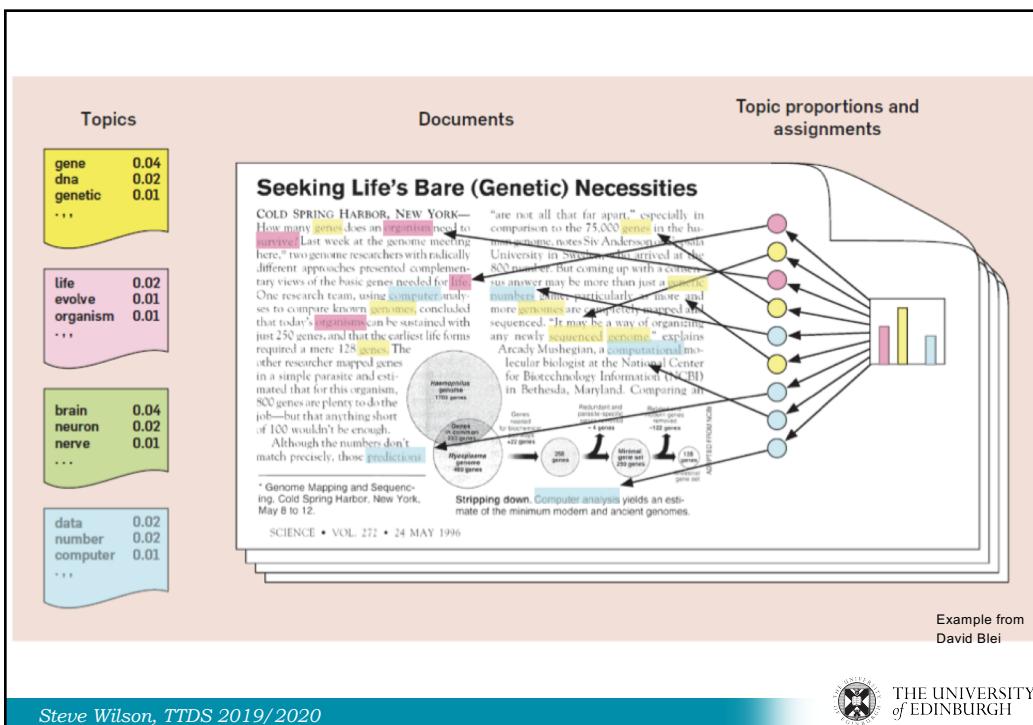
35

Steve Wilson, TTDS 2019/2020

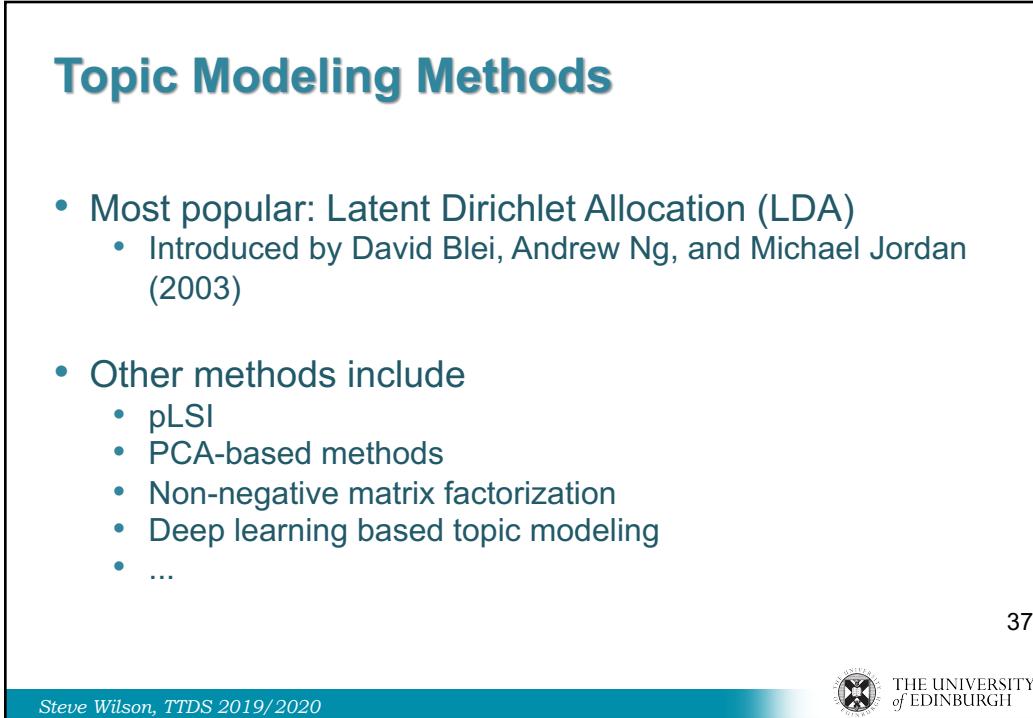


35

16



36



37

37

Topic Modeling Methods

Next Lecture

- 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
 - ...

38

Steve Wilson, TTDS 2019/2020



THE UNIVERSITY
of EDINBURGH

38

Reading

- [Manning: IR book section 13.5](#)

Steve Wilson, TTDS 2019/2020



THE UNIVERSITY
of EDINBURGH

39

18