



Healthcare Text Analytics: Analysing Free-text Health Data

Goran Nenadic^{1,2,3}

¹School of Computer Science, University of Manchester ² Health eResearch Centre ,The Farr Institute of Health Informatics Research ³Manchester Institute of Biotechnology







Aims of this tutorial

- Understand the needs, opportunities and challenges in the extraction and integration of information from free-text health data sources
- Discuss the main steps in text mining
- Show-case possible application areas
- Discuss the state of the art
- Provide some pointers for further engagement





Plan

- 1. Very brief overview of text mining
- 2. Examples of processing healthcare free text
 - Clinical notes
 - Patient generated data
- 3. Hands on
 - Extract dosage information from free text
 - Extract symptoms, procedures, anatomical locations
 - Get all patient with hypertension
- 4. State of the art
- 5. Join the community





Group reflection: Word vs. Excel

VS

- What is easier to create?
- What is easier to process?

旗	1 2 2 2 1	
-		14
	A new special state in starts descenting suppress (c) for every state, while will be appressive of the starts of the start of the starts of the start of the starts of the starts of the starts of the starts of the start of the starts of the start of the starts of the start of the starts of the start of the starts of the start of the start of the start of the starts of the start of the start of the start of the start of the start of the starts of the start of the start of the start of the start of the start of the start of the starts of the start of the start of the start of the start of the start of the start of the start of the start of the start of the start of the start of the start of the start of the start of the start of the start of the start of the start of the start	
	Index is a data series we assume a factory in the factor of the factor of the series o	
	The design over the period of the state of the local distribution of the state of t	
	Dass er som somer fräk ånden frag om som i ben være stilt frag ogsånn signer, ga forfan for Dasser som somer i somer som somer er som er som er som som er som performerer for an er som som er somerer at som for an er som er som er som er som er som er som som er som er som er som er som er som er som e	

					and galaxies							(A)
	Home	Insert	Page Layo	ut Forn	iulas D	ata Revi	ew View	v Devel	oper			· · ·
B	à	Arial	14		G	eneral *	E Cond	itional Form	atting *	** Insert *	2 7	1
Paste	-0	BIU-	A A			· % ·	Forma	at as Table *	j.	Delete +	Sort	& Find &
-	3	H - 0 - 1	A -	课课》		se .00 €.€	Cell S	tyles -		Format -	2" Filte	- Select -
Clipboa	rd 🐨	Font	G	Alignmen	t G	Number G		Styles		Cells	Edi	ting
	A1	-	6	fx Real G	DP by Met	tropolitan A	rea (milli	ons of chai	ned 2005 (dollars)		
1	A	В	C	D	E	F	G	н	E	J	К	L
1 Re	eal G	DP by N	letropo	litan A	rea (mil	lions of	chaine	ed 2005	dollar	S)		
2												
3 Bur	reau of	Economic Ar	alvsis									
4 All	industr	v total	- C									
5												
6	Fips	Area	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010
51 470	20	Victoria, T.	3718	3864	3895	4197	4262	4433	4484	4341	4010	4167
52 472	220	Vineland-N	4048	4047	4179	4299	4408	4446	4417	4466	4369	4383
53 472	260	Virginia Be	60334	62105	64426	66119	68326	70367	71722	72014	71277	71557
54 473	300	Visalia-Por	7946	8255	8702	9461	10025	10018	10597	10295	9903	9997
55 473	380	Waco, TX	5958	6091	6343	6658	6856	7013	7270	7397	7426	7739
56 475	680	Warner Ro	3930	4117	4197	4279	4468	4563	4664	4599	4698	4810
57 479	000	Washingto	294656	304317	316043	333191	348752	354687	360390	369763	369771	383073
58 479	940	Waterloo-C	5575	5890	5874	6560	6695	6678	6937	6914	6739	6880
59 481	140	Wausau, V	4842	4779	4980	5087	5231	5277	5364	5160	5028	5047
60 483	300	Wenatchee	2693	2855	2983	3040	3101	3229	3197	3260	3283	3265
61 485	540	Wheeling,	4066	4159	4268	4346	4335	4315	4275	4301	4476	4572
62 486	520	Wichita, K	23313	23029	22595	22188	22393	24735	26185	25515	23550	23453
63 486	660	Wichita Fa	5120	5225	5146	5148	4855	5027	5165	5268	5366	5294
64 487	00	Williamspo	3236	3179	3210	3276	3303	3293	3283	3203	3172	3419
65 489	000	Wilmington	10717	10219	10330	10683	11418	11558	12290	12256	12070	12309
66 490	020	Wincheste	3735	3772	3999	4112	4428	4569	4496	4320	4304	4476
867 491	180	Winston-S	19843	18629	19464	19860	20559	20733	20986	20253	19616	20090
868 493	340	Worcester	24879	24951	25641	25671	25745	25958	26118	26472	25168	26262
869 494	120	Yakima, W	5599	5765	5923	5979	6286	6442	6524	6641	6607	6490
370 496	520	York-Hano	11482	11588	12265	12586	13261	13215	13508	13587	12935	13388
871 496	660	Youngstow	16380	16841	16776	16863	16961	16641	16297	15542	14272	14807
372 497	00	Yuba City,	3390	3654	3887	3965	3970	4118	4080	4194	4246	4169
873 497	40	Yuma, AZ	3425	3847	3856	4211	4365	4564	4678	4475	4463	4417
874 Le	egen	a / Foot	notes:									
75 NA	AICS In	dustry detail	is based or	the 2002	North Ame	rican Indust	ry Classific	ation Syste	m (NAICS).		
376 (D)	Not sh	own in order	to avoid the	disclosure	of confide	ntial informa	ation; estin	nates are in	cluded in h	igher level I	totals.	
877 (L)	Less th	nan \$500,000	in nominal	or real GD	P by metro	politan area						
378 (NA	A) Not a	wailable.										
379 (NA	() Not n	neaningful.										
880 NC	DTE: O	n September	29, 2011,	statistics o	per capita	real GDP v	rere update	d to incorp	orate Cens	us Bureau i	midyear pop	ulation
881 La	st upda	ited: Septem	ber 29, 201	1								
4 4 4	H Sh	eet0 PJ			_			0				
		and the second se							Loren 12		0	6

Natural language data is everywhere

- Internet
 - Newswire, reports, blogs, etc.
 - Facebook: 510,000 comments posted every minute (with 293,000 status updates and 136,000 photos)
 - **Twitter**: 6,000 per second (i.e. 200 billion per year)
- Scientific literature
 - ~1 million scientific articles per year in biomedicine only
- Internal company reports
 - One legal company has a repository of 1 billion documents (including different versions)
- Clinical records
 - The Christie Hospital: 135k patients





Actionable information



Structured variables and coded information

e.g. lab-tests, codes, prescriptions

Unstructured or semi-structured information

e.g. clinical narrative, letters, social media, etc.





Group reflection: what is in free-text?

 Do you have a clinical example of information that is mostly or exclusively available in free-text in your domain?



Natural Language Systems

- Started around 1950s
- Various terminology
 - natural language processing (NLP)
 - text mining

(now widely used as an umbrella for large variety of NLP techniques to denote all approaches to retrieve, extract and analyse text`)

- text analytics
- computational linguistics
- (human) language technology
- Many sub-areas





What does it take to be a data scientist

High Level Skills						
Bachelors	Masters	PHD				
ML_NEURAL_NETS_SVM	ML_NEURAL_NETS_SVM	ML_NEURAL_NETS_SVM				
OPTIMIZATION	OPTIMIZATION	TEXT_ANALYTICS				
DATA_WRANGLING	CLUSTERING	OPTIMIZATION				
CLUSTERING	DATA_WRANGLING	LINEAR_REGRESSION				
DATA_VISUALIZATION	TIME_SERIES	CLUSTERING				
TIME_SERIES	LINEAR_REGRESSION	TIME_SERIES				
TEXT_ANALYTICS	TEXT_ANALYTICS	DESIGN_OF_EXPERIMENTS_AB				
LINEAR_REGRESSION	DATA_VISUALIZATION	DATA_VISUALIZATION				
DESIGN_OF_EXPERIMENTS_AB	DESIGN_OF_EXPERIMENTS_AB	DATA_WRANGLING				
BASIC_STATISTICS	BASIC_STATISTICS	LOGISTIC_GLM				



http://blogs.sas.com/content/text-mining/2015/03/27/whats-it-take-to-be-a-data-scienti

Text Mining in one slide

- On May 2, 2014 the then Indian Prime Minister
 Manmohan Singh revealed: 'General Musharraf and I had nearly reached an agreement, a nonterritorial solution to all problems..."
- Following in suite of the electric vehicle manufacturer, General Motors has revealed its plans of growing its energy storage market.





Typical text mining framework



- Information Retrieval (IR)
 - find relevant documents [unstructured form]

Information Extraction (IE)

- find detailed information [(semi-)structured form]
- Data Mining (DM)
 - find associations, build networks, make predictions

Information retrieval (IR)

- Searching for relevant documents
 - Find similarity between user query
 - Rank documents
- Result of IR is a set of relevant documents
 - filtering huge collections based on a query
 - no fine-grained information, just whole documents
 - users would need to read and analyse these documents on their own





Information extraction (IE)

- Extract information i.e. facts from text
- Identify instances of pre-defined *entities* (dates, names of people, locations, etc.) and relations between them
- Fill in database-like tables with "facts"

Slot	Information
Date	7/10/96 (today)
Location	SanSalvador
Victim injured	policeman
Victim attacked	guards
Perpetrator	urban guerrillas

San Salvador, 7/10/96

It has been officially reported that a policeman was wounded today when urban guerrillas attacked the guards at a power substation located downtown San Salvador.



Main problems: ambiguity

- Different meanings/senses of words
 - e.g. Apple (the company) or apple (the fruit)
 - e.g. *Toyota* can be a *car* or a *company*
 - e.g. acronyms have different meanings in different contexts

United States Army	Union of South Africa		
United States of America	Union Street Athletics		
Ulhasnagar Sindhi Association	Unionville-Sebewaing Area		
Ultimate in Suspense and Action	Unique Settable Attributes		
Unconditional Self-Acceptance	Unit Self Assessment		
Unconventional Stellar Aspect	United Scenic Artists		
Under Secretary of the Army	University of South Alabama		
Underground Service Alert	University of South Australia		
Underground Sewer Adapter	Unix System Admin		
Underwriting Service Assistant	Unstable Angina		
Unicycling Society of America	Unusually Sensitive Area		

USA =

Main problems: variability

Numerous ways to express the same thing

Chronic renal impairment Chronic kidney disease CKD kidney failure Chronic renal failure CRF Chronic renal failure syndrome Renal insufficiency eGFR 44







Basic text pre-processing

- Basic \neq trivial
- Typically includes:
 - Identify words (tokenisation)
 - oxycodone 5-10 mg p.o. q.4 h. as needed for pain.
 - Normalise <u>word</u> forms
 - e.g. plural forms, or past tense verbs
 - Segment sentences
 - punctuation not always used!







Main steps in text processing



Finding entities/concepts

- Identification of <u>entities</u> and <u>terms</u> of interest
 - find strings in text that denote specific entities or concepts (named-entity recognition = NER)
 - e.g. persons, jobs, organisations, dates, time, locations, symptoms, disorders, medications, treatments, ...
 - designed for each *class* of interest





NER example

🛎 Entity determination.

File Entities

Add entity **Remove entity** Enter text to search: **Contract entities** AutoContraction **Reset highlight** Search for selected Add as entity id person_name comment id organisation_n... 34 Tolva Shponar pliot:Belarus ٠ 16 Ihangvi REPORT DATE : 9 October 2008 | provided to CIA by Pakistani Criminal 36 Raph Plotnitsk ... pllot; Ukra. 14 CCTV hvestigation Unit, Karachi Division | NOTE: Surveillance report on the 37 Pepik Malakho., pilot Ukrai. 32 Crime Suppres. activities of Maulana Hag Bukhari, suspected to be a top leader within the 38 Mykola Khitovo pilot Ukrai. 47 Maulana Hag Karachi faction of Lashkar-e-Jhangvi Bukhari is frequently accompanied by 39 Grigor Jalovskaja pilot Ukraj. dram Baera, who acts as a driver and bodyguard . Additional information 50 Bukhari 42 Arkadi Borodi... Ukraine;b.. provided by police informants 55 Katchi Abadis 51 Akram Basra 23 July 2008 - A delivery was made to a house in Lyari Town (a 13 CID 4 . 4 constituent town of Karachi) in a house believed to be used by Bukhar id location city country ic MISC The delivery was made by a two men in street clothing (as opposed to a A 19 Ilvushin IL-76 29 uniform) who arrived in a white van , license **(1) and the set of** , with single blue Sri Lanka 45 Karachi Division 35 Belarus. stripe on each side The delivery consisted of three medium boxes (requiring two hands to move) and a small box (handsized) The boxes 40 Ukraine appeared to be heavy. One large box was square, the other rectangular 46 Surveillance The small box was rectangular. It is unknown if But har was home at the 49 Lashkar-e-lha. time of the delivery. 52 Lvari Town 8 August 2008 - An unknown man visited the Lvari Town house where + 48 Karachi turnate is believed to stay. He arrived at 1615, and was let into the house ЪĨ 1 • immediately. Loud voices could be heard for a few moments, then they id tool date id. subsided About fifteen minutes later a silver Mercedes left the rear of the 43 February 11 * 4 Sixty National . house with what appeared to be three occupants. Due to the tinted glass it 57 23 July 2008 6 Computerized. was impossible to identify the occupants of the vehicle. The house was 58 8 August 2008 61 LHR 6354 surveilled for the next five hours , however no one came or went 59 16 September. 16 September 2008 - Bukhari was followed from his Lyan Town house to 60 23 September. an apartment about 1.5 kilometers southeast. He entered the building and staved there for about two hours . When he left he returned to his Lvari Town house at which time the observer believed he 'd been discovered and 4 . 4 left the area Id event comment 23 September 2008 - Bukharl and probably Basra are reported to have visited a house in the Katchi Abadis Old Settlement, on 835T Longhi Street The informant, a vendor with business in the area, was passing by and noticed Bukhar and one other person enter the building at about 1430 The informant knew Bukhar by sight , but had never met him 0 entity found and highlighted. 4 . Entity definition

12161

K

.

Move entity

Remove selections

comment

comment

cargo plane

comment

Mercedes

Two kinds of NE approaches

Knowledge Engineering

- developed by experienced language engineers
- make use of human intuition
- requires only small amount of training data
- development could be very time consuming
- some changes may be hard to accommodate

Learning Systems

- use statistics or other machine learning
- developers do not need language expertise
- requires large amounts of annotated training data
- some changes may require reannotation of the entire training corpus





NER methods

- dictionary-based (gazetteers)
 - matches items from a **lexicon** against text

rule-based

- hand-written regular expressions
- e.g. <title> <capitalised word>+ → personal name
- e.g. <person>Mr John Smith Jr.</person>

machine learning

- often supervised models to identify specific entities
- hybrid any combination of the above

I. Spasic



T. Liptrot

Finding relations between entities

- Extract specific facts, relations and events by linking entities
 - use of templates, regular expressions, grammars
 PERSON> is appointed as a <JOB> of <COMPANY> ORUG> for <SYMPTOM>
 - typically designed around important verbs
 - e.g. attacked, bought, appointed, merged, acquired, etc.
 - also, machine learning approaches





Finding relations between entities

Manual annotation of a training set



Quality metrics

- The result is compared with a manual "gold standard"
 - How consistent the humans are/can be?
 - 100% is impossible even for human annotators.
- Typical measures:
 - Sensitivity (true positive rate, recall, probability of detection) proportion of positives that are correctly identified as such
 - Specificity (true negative rate) proportion of negatives that are correctly identified as such
 - Precision (positive predictive value, PPV) proportion of retrieved elements that are positives
 - F-measure (harmonic mean of sensitivity and PPV) combines sensitivity (recall) and precision





Measuring agreement (Kappa)

• Measures the agreement between two annotators who each classify *N* items into *C* mutually exclusive categories

$$\kappa = \frac{\Pr(a) - \Pr(e)}{1 - \Pr(e)},$$

- Pr(*a*) = relative observed agreement among annotators
- Pr(e) = the hypothetical probability of chance agreement
 - typically using the observed data to calculate the probabilities of each observer randomly saying each category
- If the annotators are in complete agreement then $\kappa = 1$.
- κ = 0 if there is no agreement among the annotators other than what would be expected by chance (as defined by Pr(e))





Mining healthcare free-text







Routinely collected narrative data

 Clinician notes, letters are part of routinely collected data

ALLERGIES: To IV contrast, penicillin

- People communicate most efficiently in narrative text
- Not everything can be captured in pull-down menus or check-lists
- Not everything out bo out bo

cerebrovascular accident.

OTHER DIAGNOSES PERIPHERAL VASCULAR DISEASE , HYPERTENSION.



experts





Large-scale data is available

Clinical Record Interactive Search (CRIS) system

- SLaM CRIS
- De-identified EPR database
- Structured and free text data



UK-CRIS network

 10 Mental Health NHS Trusts (total 2M patients)





NHS National Institute for Health Research

Text Mining Clinical Narrative

A 43-year-old woman was diagnosed with type 2 diabetes mellitus by her family physician 3 months before this presentation. Her initial blood glucose was 340 mg/dL. Glyburide 2.5 mg once daily was prescribed. Since then, selfmonitoring of blood glucose (SMBG) showed blood glucose levels of 250-270 mg/dL. She was referred to an endocrinologist for further evaluation.

On examination, she was normotensive and not acutely ill. Her body mass index (BMI) was 18.7 kg/m2 following a recent 10 lb weight loss. Her thyroid was symmetrically enlarged and ankle reflexes absent. Her blood glucose was 272 mg/dL, and her hemoglobin A1c (HbA1c) was 10.3%. A lipid profile showed a total cholesterol of 261 mg/dL, triglyceride level of 321 mg/dL, HDL level of 48 mg/dL, and an LDL of 150 mg/dL. Thyroid function was normal. Urinanalysis showed trace ketones.

She adhered to a regular exercise program and vitamin regimen, smoked 2 packs of cigarettes daily for the past 25 years, and limited her alcohol intake to 1 drink daily. Her mother's brother was diabetic. A 43-year-old woman was diagnosed with type 2 diabetes mellitus by her family physician 3 months before this presentation. Her initial blood glucose was 340 mg/dL. Glyburide 2.5 mg once daily was prescribed. Since then, self-monitoring of blood glucose (SMBG) showed blood glucose levels of 250-270 mg/dL. She was referred to an endocrinologist for further evaluation.

On examination, she was normotensive and not acutely ill. Her body mass index (BMI) was 18.7 kg/m2 following a recent 10 lb weight loss. Her thyroid was symmetrically enlarged and ankle reflexes absent. Her blood glucose was 272 mg/dL, and her hemoglobin Alc (HbAlc) was 10.3%. A lipid profile showed a total cholesterol of 261 mg/dL, triglyceride level of 321 mg/dL, HDL level of 48 mg/dL, and an LDL of 150 mg/dL. Thyroid function was normal. Urinanalysis showed trace ketones.

She adhered to a regular exercise program and vitamin regimen, smoked 2 packs of cigarettes daily for the past 25 years, and limited her alcohol intake to 1 drink daily. Her mother's brother was diabetic.





Text Mining Clinical Notes

DisorderCEMtext:diabetes mellituscode:73211009subject:patientrelative temporal context:3 months agonegation indicator:not negated	A 43-year-old woman was diagnosed with type 2 diabetes mellitus by her family physician 3 months before this presentation. Her initial blood glucose was 340 mg/dL. Glyburide 2.5 mg once daily was prescribed.
Medication CEMtext:Glyburidecode:315989subject:patientfrequency:once dailynegation indicator:not negatedstrength:2.5 mg	
TobaccoUseCEMtext:smokingcode:365981007subject:patientrelative temporal context:25 yearsnegation indicator:not negated	She adhered to a regular exercise program and vitamin regimen, smoked 2 packs of cigarettes daily for the past 25 years
Disorder text:CEMtext:diabetes mellituscode:73211009subject:family memberrelative temporal context:not negated	Her mother's brother

Strategic Hustiti 17 Advanced Research Projects (SHAISP) Program

dependent of the United States and American

statute second statute



Text Mining Clinical Notes

Transform free-text into structured data: extract clinical variables and their values

Diseases, problems, symptoms

- acute dizzyness, nausea
- Anatomy
 - upper eyelid, tongue, throat
- Drugs/medications
 - steroids, desloratadine, zoladex
- Adverse drug events
 - differentiate from symptoms
- Diagnostic/therapeutic procedures
 - hysterectomy
- Behaviour/quality of life
 - anxiety, problems at work, poor sleep





Mining clinical narrative

Examples with real data

- Extracting baseline patient data
- Extracting dosages from free text prescriptions – finding drug regiments
- Mining patient journeys
- De-identification of clinical narrative
- Symptom severity in psychiatric evaluation notes
- What do EHR talk about?

Extracting baseline patient data

- 6 elements:
 - primary disease site
 - histology
 - stage
 - performance status
 - comorbidities
 - responsible consultant
 - date of diagnosis & death
- 135,000 sets of notes manual work would take
 80 years





Approaches – rule-based

• Create *rules* that are matched against text

T[Xx0-4isd][abcdismi]* N[Xx0-4is][abcdismi]* M[01Xx]

> "...scan showed a stage of T3N4M0..." "Current diagnosis showed TxN0M0..."

Results

class	method	precision	recall	f_measure	extras	
site	abstractors*	0.83	NA	NA	NA	
	rules	0.84	0.61	0.71	0.01	
	ml	0.96	0.53	0.68	0.00	
	ml + rules	0.88	0.72	0.79	0.01	
histology	abstractors*	0.92	NA	NA	NA	
	rules	0.89	0.71	0.79	0.02	
	ml	0.83	0.32	0.46	0.09	
	ml + rules	0.84	0.82	0.83	0.10	
t stage	rules	0.84	0.70	0.76	0.05	
	ml	0.90	0.75	0.82	0.09	
	ml + rules	0.91	0.85	0.88	0.12	
n stage	rules	0.90	0.55	0.68	0.00	
	ml	0.94	0.62	0.75	0.21	
	ml + rules	0.91	0.81	0.86	0.21	
m stage	rules	1.00	0.36	0.53	0.00	
	ml	1.00	0.41	0.59	0.19	
	ml + rules	1.00	0.55	0.71	0.19	
ps	rules	1.00	0.87	0.93	0.07	
	ml	0.97	0.76	0.85	0.23	
	ml + rules	0.97	0.92	0.95	0.25	
responsible consultant	ml	0.85	0.75	0.79	0.03	
	abstractors*	0.84	NA	NA	NA	The Christie



Mining clinical narrative

Examples with real data

- Extracting baseline patient data
- Extracting dosages from free text prescriptions – finding drug regiments
- Mining patient journeys
- De-identification of clinical narrative
- Symptom severity in psychiatric evaluation notes
- What do EHR talk about?

Medication prescriptions

5mg prednisolone o.d. for 2/52	
Janes -	

	studyno	start	stop	
onarug	innéénni.	11/01/7002	77/05/1002	
	00000000	17/24/2000	14/07/1001	0
10	organities -	14,101,72000	20/07/1002	1
	hoodood.	10/02/2000	11/08/1001	1
(8)	1100000	85/05/2000	11/10/0003	
14	00000001	12/10/2002	22/20/1002	1
1	dollanns	12/16/2002	37,702,71003	1
140	0.000000	17/00,22000	13/07/0003	- A.
	000000011112	\$9,70772005	812/08/3884	.0
1882	0000008	0276672004	20/11/2005	0
128	0.000000	19/11/2004	85,097/0002	A.
- 11	dom/state	81/01/2002	34/07/3882	
100	0000001	14/07/2005	12/10/2002.	.0.
124	0000002	17/10/2000	11/07/2004	4
24	andorum ?	83/07/2004	21,/02,/3.00k	8
(34.)	0000001	11/02/2000	20/13/2009	0
1.9	0000008	10/01/2000	18/01/0804	1
188.		11,702,72004	14/02/3304	
188	0000001	15/03/2004	13/04/1904	(Q)
0.0	0000000	11/66/2004	11/02/5004	
144		18,786,72006	342/00/3388	
3.1	0006001	14/02/2000	10/21/2905	0
23	0000000	10/01/2001	\$4,707,70008	0.
-24-		34/82/2004	30/22/3008	
38	0000000	80/01/2008	\$4/84/1903	0
28	0000000	1*/00/2000	12/09/0000	0.
:27:	0	17785/2000	28/02/2885	8
34	00000000	19/00/2000	10/12/2008	0
28	0006008	10/01/2000	14/07/2003	1
10.	0	84/07/2000	\$870772202	
	0088008	14/07/2805	18/09/2003	10.
11	0.000000	15/05/2000	04/03/2005	0

Raw Data

Data prepared for analysis

Medication prescriptions

- Dose, frequency, duration, intervals
 - take 2 tablets 4 times a day for 2/52
- But often complex prescriptions
 one or two to be taken every 4 to 6 hours
 - 10mg to be taken weekly
 - a half to one tablet to 2 three times a day
 when required



Karystianis et al, BMC Med Inform Decis Mak. 2016;16:18. doi: 10.1186/s12911-016-0255-x.



Modelling prescription data

Prescription		dose		freq		interval	
		max	min	max	min	max	unit
take 2 tablets 4 times a day	2	2	4	4	1	1	tablet
a half to one tablet to 2 three times a day <u>when</u> <u>required</u>	0.5	1	0	3	1	1	tablet
10mg to be taken <u>weekly</u>	10	10	1	1	7	7	mg
2 with <u>each meal</u>	2	2	3	3	1	?	?
1 to 3 every day	1	3	1	1	1	1	?
one or two to be taken every 4 to 6 hours	1	2	4	6	1	1	?

Variability in prescriptions

- 56,000 most common free text instructions from **CPRD** records
 - CPRD = a UK Anonymised Primary Care EHR Research Database, with over 11M patients

- at least 1 in 4 has inherent variability

 i.e. a choice in taking medication specified by different minimum and maximum doses, duration or frequency.



Karystianis et al, BMC Med Inform Decis Mak. 2016;16:18. doi: 10.1186/s12911-016-0255-x.



Mining clinical narratives

Examples with real data

- Extracting baseline patient data
- Extracting dosages from free text prescriptions – finding drug regiments
- Mining patient journeys
- De-identification of clinical narrative
- Symptom severity in psychiatric evaluation notes
- What do EHR talk about?

Mining Patient Journeys

- Can text mining techniques be used to <u>reconstruct patient journeys/pathways</u> from unstructured clinical narratives?
- Data: longitudinal narratives
 - Clinical narratives: internal notes and correspondence letters
 - Study participants (n = 27): adult survivors of childhood central nervous cancer (medulloblastoma)

The Christie M/S

NHS Foundation Trust

"Comparing the narrative experiences of patients with medulloblastoma with factors identified from their hospital records"

Mining Patient Journeys

- Healthcare concept extraction
 - problems (e.g., headache, brain cancer),
 - *treatments* (e.g., chemotherapy, radiotherapy)
 - <u>tests</u> (e.g., MRI scan, blood test)
 - <u>health-related quality of life</u> (e.g., physical, emotional functioning)

Temporal Information Extraction

- Recognition and normalisation of temporal expressions
- Chronological ordering of concepts

Clinical temporal expressions

ADMISSION DATE: 2011-02-06 DISCHARGE DATE: 2011-02-08



M. Filannino

Mining Patient Journeys

Clinical Dautoisand Home

Satting +



Mining Patient Journeys

The second se	Occurry Divergences		
Parlamente (1993) Annos Josephan Santar 1993 Ange M	Patient ID MB001 Celegory problem Type	Semantic Type Neoplassic Interest Semantic Groop Disorders Temporal 3/5/1990 Intermution	→ ₩ 0 + Pimr
reactive and Carly HUAP	Comment extendences Document ID Heference 2 MED01-229 3/5/1989 Lorom issuer dotor sit amer, consected, at amer operate quan camue su. Prov Maura consequal felia su taplen posue Hec thomas ums preture. Present sus previae ecte qua libero volutpat volutpat	Nate Type Climic Hoter In schpielong ett. Proin riturneuw ex longth, Feribust tellus at suito mattile portfilitet re-autor. Ut lobortie leo pus faile effectur, emod ribh at alliquet venerabits, triteger st. Integer au factinis anto.	2010 20
		000	

Patient TimeLanes @ cTAKES



https://sourceforge.net/projects/ohnlp/files/cTAKES/

Mining clinical narratives

Examples with real data

- Extracting baseline patient data
- Extracting dosages from free text prescriptions – finding drug regiments
- Mining patient journeys
- De-identification of clinical narrative
- Symptom severity in psychiatric evaluation notes
- What do EHR talk about?

De-identification of narrative

• Access to data?



HIST	TORY OF PRESENT ILLNES	SS:	The patie	ent is	a 77-y	ear-old-woman	
with	long standing hypertension v	who	presente	ed as	a wall	k-in to me at	
the	the Oak Valley Health Center on July 9th . Recently had been						
start	started q.o.d. on Clonidine since May 5th to taper off of the drug.						
Was	Was told to start Zestril 20 mg. q.d. again. The patient was sent to						
the Smith Cardiac Unit for direct admission for cardioversion and							
antic	anticoagulation, with the Cardiologist, Dr. Pearson to follow.						





Automated de-identification

- Motivation
 - Enable data access to researcher
 - Ethical and legal requirements
 - Time and cost vs. manual labour
- Personal identifiable Information (PII)
 - Name, Date of Birth, Contact, Address, Phone number, ...
 - NHS number,
 - Profession, family members, etc.
- E.g. clinical letters at the Christie hospital
 - Accuracy 96.05% comparable to human benchmark



Mining clinical narrative

Examples with real data

- Extracting baseline patient data
- Extracting dosages from free text prescriptions – finding drug regiments
- Mining patient journeys
- De-identification of clinical narrative
- Symptom severity in psychiatric evaluation notes
- What do EHR talk about?

Symptom severity in psychiatric evaluation notes

- Analyse patient initial evaluation notes to identify symptom severity (absent, mild, moderate, severe)
 - Does a patient requires immediate medical attention or hospitalisation?
 - Based on the Research Domain Criteria (RDoC) framework
 - Focus on: positive valence, negative valence, cognitive, social processes, arousal and regulatory system
- Input data is a mix of
 - semi-structured questions (e.g. "History of Drug Use: Yes") and
 - free-text narrative (e.g. *"treatment included x2 IOP treatment, IOP for alcohol dependence"*)











Symptom severity in psychiatric evaluation notes

- Methods:
 - Knowledge-based approach (rules)
 - Identify key-phrases (e.g. substance use, manic episodes, treatments, etc.)
 - Data-driven approach (deep learning)



- Current results
 - absent: 82%; mild: 93%; moderate: 70%; severe: 77%











Mining clinical narrative

Examples with real data

- Extracting baseline patient data
- Extracting dosages from free text prescriptions – finding drug regiments
- Mining patient journeys
- De-identification of clinical narrative
- Symptom severity in psychiatric evaluation notes
- What do EHR talk about?

What do EHR talk about?













What patients are talking about?

















Healthcare Text Analytics: Demos

George Demetriou¹

¹School of Computer Science, University of Manchester





Plan

- 1. Very brief overview of text mining
- 2. Examples of processing healthcare free text
 - Clinical notes
 - Patient generated data
- 3. Hands on
 - Extract dosage information from free text
 - Extract symptoms, procedures, anatomical locations
 - Get all patient with hypertension
- 4. State of the art
- 5. Join the community

Task: Extract Dosage Information

Variability in text:

- "25 mg" -> Numeric_value (integer) + Unit
 - -> Numeric_value (decimal) + Unit
- "20-25 mgs"

• "15.5 kg"

- -> Numeric_value (range) + Unit
- "6 hours" -> Time: different feature of dosage!
- What if we want to match

e.g. 'tblspoons' or 'tbsp' ?