

# Helping the medicine go down? Repair and adherence in patient-clinician dialogues

Matthew Purver, Christine Howes, Rose McCabe

Interaction, Media and Communication  
School of Electronic Engineering and Computer Science  
Queen Mary, University of London  
[m.purver@qmul.ac.uk](mailto:m.purver@qmul.ac.uk)



## Interaction, Media & Communication

- Multi-disciplinary research group
- Human-computer interaction
- (Computer-mediated) human-human interaction
- Non-verbal and verbal communication
  - Gesture and posture in interaction
  - Incremental NLG and NLU for dialogue
  - Communication on social media
- EPSRC proof of concept project
  - Joint work with Barts & the London SMD
  - Application of NLP to doctor-patient dialogue

# Outline

- 1 Background
  - Doctor-patient communication and adherence
  - Repair
- 2 Repair in patient-doctor communication
  - Method
  - Results
- 3 Automatic detection of repair
  - Approach & previous work
  - Results
  - Next steps
- 4 Automatic prediction of adherence
  - Approach & previous work
  - Results
  - Next steps

# Outline

- 1 Background
  - Doctor-patient communication and adherence
  - Repair
- 2 Repair in patient-doctor communication
  - Method
  - Results
- 3 Automatic detection of repair
  - Approach & previous work
  - Results
  - Next steps
- 4 Automatic prediction of adherence
  - Approach & previous work
  - Results
  - Next steps

# Doctor-patient communication

- Non-specific effects account for 60% of variance in patient outcome in clinical trials (Walach et al, 2005)
- One locus of non-specific effects is doctor-patient communication
- Shared understanding important (Mead et al, 2000)
  - How the patient understands the doctor
  - Common understanding of goals and implementation of treatment
- Associated with patient outcomes (Ong et al, 1995)
  - Patient satisfaction
  - Treatment adherence

# Doctor-patient communication in schizophrenia

- Non-adherence to treatment a significant problem with schizophrenia
  - About half of patients are non-adherent in the year after discharge from hospital (Weiden & Olfson, 1995)
  - Risk of relapse 3.7 times higher (Fenton et al, 1997)
- Shared understanding challenging in schizophrenia (McCabe et al, 2002)
  - Problems in communicating and understanding
  - Patients may not accept explanations of delusions
  - Psychiatrists focus on treatment may not be in line with patients perspective
- Recent research suggests *repair* associated with adherence
  - (McCabe et al., in prep.)

## Repair is ...

- A dialogue phenomenon – used by speakers to:
  - formulate understanding of one's own talk
  - clarify understanding of other's talk
  - address misunderstanding of own and other's talk
- Pervasive in dialogue (e.g. Schegloff, 1992)
  - identifying and resolving (potential) misunderstandings
  - driving the process of grounding (Clark, 1996)
  - collaboration to achieve shared understanding
- Important in dialogue systems
  - signalling (potential) misunderstandings
  - understanding corrections and confirmations

## Types of repair

- Repair classified according to position (which turn), initiator (self or other) and repairer (self or other)



## Types of repair

- Repair classified according to position (which turn), initiator (self or other) and repairer (self or other)

### Example - Turn positions

A:	How are you?	← <i>position 1</i>
B:	Fine thanks	← <i>position 2</i>
A:	That's good	← <i>position 3</i>

## Types of repair

- Repair classified according to position (which turn), initiator (self or other) and repairer (self or other)
  - P1SISR - Signal problem and resolve in own turn

### Example - P1SISR - Articulation

Dr: You probably have seen so many psychiatrists **o o over the years**

### Example - P1SISR - Formulation

Dr: **Did you feel that did you despair so much that** you wondered if you could carry on

### Example - P1SISR - Transition Space

P: Where I go to do **some printing lino printing**

## Types of repair

- Repair classified according to position (which turn), initiator (self or other) and repairer (self or other)
  - P3SISR - Signal problem and resolve in subsequent own turn

### Example - P3SISR

Dr: **Clorazil** or

P: Yeah

Dr: **Clozapine** yes

## Types of repair

- Repair classified according to position (which turn), initiator (self or other) and repairer (self or other)
  - P2OIOR - Other person signals and resolves problem in next turn

### Example - P2OIOR

Dr: Rather than **the diazepam** which I don't think is going to do you any good

P: **the valium**

## Types of repair

- Repair classified according to position (which turn), initiator (self or other) and repairer (self or other)
  - P2NTRI - Other person signals a problem for the original speaker to resolve
  - P3OISR - Original speaker resolves a problem signalled by the other

### Example - P2NTRI with P3OISR

Dr: Yeh, it doesn't happen in real life does it?

P: **What do you mean by real life?**

## Types of repair

- Repair classified according to position (which turn), initiator (self or other) and repairer (self or other)
  - P2NTRI - Other person signals a problem for the original speaker to resolve
  - P3OISR - Original speaker resolves a problem signalled by the other

### Example - P2NTRI with P3OISR

Dr: Yeh, it doesn't happen in real life does it?

P: **What do you mean by real life?**

Dr: **You can't - there are no messages coming from the television to people are there?**

# Outline

- 1 Background
  - Doctor-patient communication and adherence
  - Repair
- 2 Repair in patient-doctor communication
  - Method
  - Results
- 3 Automatic detection of repair
  - Approach & previous work
  - Results
  - Next steps
- 4 Automatic prediction of adherence
  - Approach & previous work
  - Results
  - Next steps

## Recent study

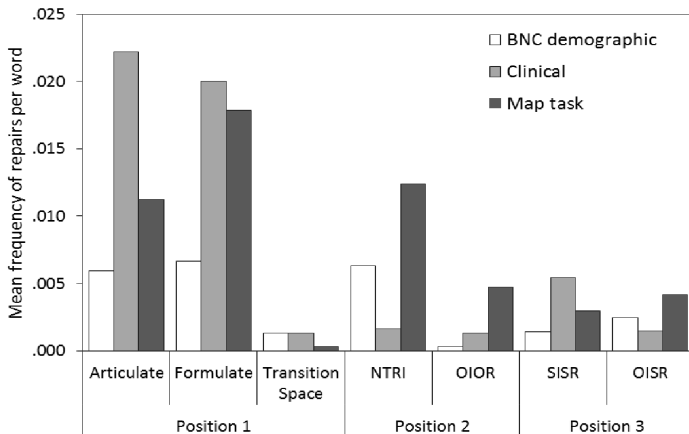
- Observational study of psychiatrist-patient consultations
  - McCabe et al. (in prep.), Howes et al., (2012)
- Assess the use of repair in negotiating shared understanding
- Test the hypothesis that more repair is associated with better treatment adherence
- Explore which types of repair are relevant for adherence
- Compare with other dialogue data (Colman & Healey, 2011)



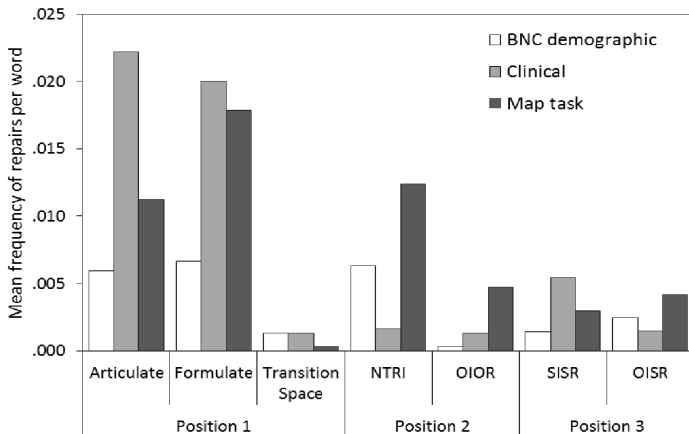
# Design

- 138 consultations audio-visually recorded
  - Consultations transcribed and repair annotated
  - Inter-annotator agreement good ( $\kappa = 0.73$ )
- Patients interviewed to assess symptoms and evaluate communication
  - 30 item Positive and Negative Syndrome Scale (PANSS)
  - Patient Experience Questionnaire (PEQ)
- Adherence assessed after 6 months (general/medication)
  - Good >75%
  - Average 25-75%
  - Poor <25%

# Comparison with other dialogue contexts

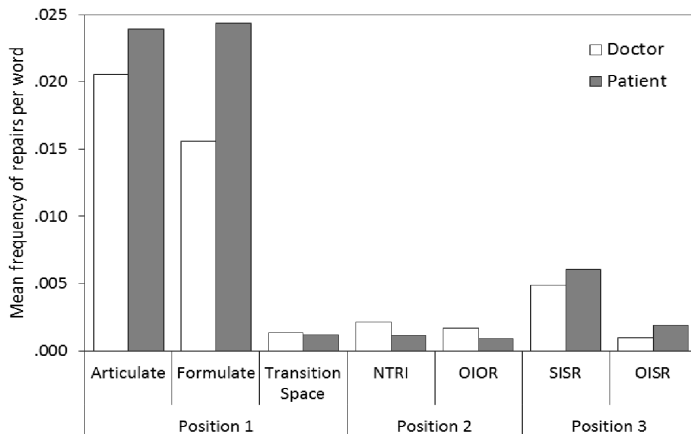


## Comparison with other dialogue contexts

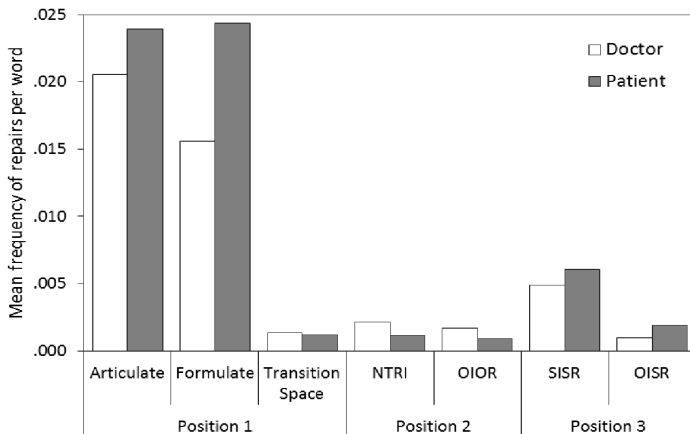


- More self-repair in clinical context
- Fewer NTRIs in clinical context

# Mean frequencies of repair types per word



## Mean frequencies of repair types per word



- Patients do more formulation repairs than doctors
- Patients produce fewer NTRIs than doctors

## Repair and adherence

- Response variable (adherence) binary; good (>75%) or not
- Adherence not associated with length of illness or symptoms
- Principal component analysis to reduce number of variables
  - 4 factors, explaining 72% of the variance:
    - Psychiatrist led clarification and patient response (31%)
    - Patient led clarification and psychiatrist response (17%)
    - Patient reformulation (14%)
    - Psychiatrist reformulation (9%)
- Regression model (mixed effects)
  - One significant association with adherence after 6 months
  - Patient led clarification: odds ratio 5.82,  $p = 0.02$
  - But 95% confidence margin wide: 1.3-25.8

## Questions

- If repair correlates with adherence:
  - can we automatically detect repair?
  - can we automatically predict adherence?

## Questions

- If repair correlates with adherence:
  - can we automatically detect repair?
  - can we automatically predict adherence?
- What does patient-led clarification measure?
  - Engagement
  - Understanding
  - Confrontation
  - Psychiatrist's communication style



## Questions

- If repair correlates with adherence:
  - can we automatically detect repair?
  - can we automatically predict adherence?
- What does patient-led clarification measure?
  - Engagement
  - Understanding
  - Confrontation
  - Psychiatrist's communication style
- What else correlates with adherence (or with other outcomes)?
  - Other dialogue phenomena
  - Content of discussion

## Examples

- Time for some quick examples ... ?

# Outline

- 1 Background
  - Doctor-patient communication and adherence
  - Repair
- 2 Repair in patient-doctor communication
  - Method
  - Results
- 3 Automatic detection of repair
  - Approach & previous work
  - Results
  - Next steps
- 4 Automatic prediction of adherence
  - Approach & previous work
  - Results
  - Next steps

## Task: detecting patient-led clarification

- Identify repair phenomena of interest
  - next turn repair initiators (NTRI)
  - position 2 repairs (P2R)
  - especially in patient talk
- Automatically annotate transcripts for relevant features
  - following e.g. (Purver et al, 2001)
- Supervised discriminative classification
  - following (Fernandez et al. 2006, Schlangen 2005)

# Dialogue act tagging

- Related to the general dialogue act tagging task
  - (label each turn with indication of function)
  - DA tagsets often include e.g. check category
- This is a sparse phenomenon
  - little attention paid in general tagging
  - (less common than question, statement etc)
- Accuracies on repair-related DA tags generally low:
  - e.g. Surendran & Levow (2006) on text
  - check 8% turns: 45% f-score
  - clarify 4% turns: 19% f-score
- Even sparser in our data: 0.8% of patient turns

# Emotions in suicide notes

- Liakata, Kim, Saha et al (2012)
- Tagging sentences into a set of categories
- Some categories common (90% of data)
  - instructions, information
  - love, hopelessness, guilt, blame
- Sparse for some categories
  - happiness, hopefulness, fear, pride, abuse . . .
- Union of set of supervised classifiers
  - Sequence classifiers (CRFs) as well as SVMs
  - 45.6% f-score overall
  - Sparse categories as low as 5% f-score

## Specific dialogue phenomena

- Existing work for similar dialogue phenomena
- Fernández, Ginzburg and Lappin (2006)
  - Classifying *non sentential utterances* (fragments)
  - Supervised classification, lexical/dialogue features
  - Accuracy very good (>90% for clarification)
  - Only a subset of our task:
    - fragments only, fragment already identified
- Schlangen (2005)
  - Finding fragments (and antecedents)
  - Fragments relatively sparse: 5% of turns
  - Finding fragments accuracy low (30-40% f-score)

## Clarification request (CR) taxonomies

- Purver, Ginzburg and Healey (2001)
- CRs are requests for the other person to provide repair of some aspect of a prior turn – subset of NTRIs
- Taxonomy of surface forms that CRs can take including
  - reprise fragments and sentences ( *“a handbag?”* )
  - wh-fragments ( *“who?”* )
  - wh-substituted reprises ( *“have I what?”* ),
  - explicit questions ( *“what do you mean by that?”* )
  - general forms ( *“eh?”*, *“pardon?”* )
- See also (Schlangen & Rodriguez, 2004)



## Clarification request (CR) taxonomies

- Purver, Ginzburg and Healey (2001)
- CRs are requests for the other person to provide repair of some aspect of a prior turn – subset of NTRIs
- Taxonomy of surface forms that CRs can take including
  - reprise fragments and sentences (*“a handbag?”*)
  - wh-fragments (*“who?”*)
  - wh-substituted reprises (*“have I what?”*),
  - explicit questions (*“what do you mean by that?”*)
  - general forms (*“eh?”*, *“pardon?”*)
- See also (Schlangen & Rodriguez, 2004)
- No automatic detection/classification
- But can provide basis for distinctive features

## Using specific lexical items

Dr: Ok you have done it before

P: **Pardon?**

Dr: If you have done it before

Dr: Presumably the ice has gone

P: **Eh?**

Dr: Presumably the ice has gone, it was quite icy this morning

Dr: Now from the psychiatric point of view because I'm not really a physical doctor that is your GP, who is your GP now

P: **What?**

Dr: Who is your GP

## Using specific sentences

P: They're not negative erm but they're positive as i eh erm  
um it's like imagining how your life will be

Dr: Ok, ok, ok so thinking about how

P: **Do you know what I'm talking about?**

Dr: What, what you want to achieve in the future what you  
want to do

Dr: Well what kind things when you see yourself and you say  
you want to go back to to where you left of how you see  
yourself

P: **I'm not with you**

Dr: How do you look at yourself as in do you see positive  
things do you see negative things

## Repeating lexical items

Dr: Yep well that is a possible side effect

P: **Side effect?**

Dr: Of the err Haliperidol

Dr: One thing that I ask you is when you were low in mood  
did you have suicidal thoughts

P: **Did I have ... ?**

Dr: Suicidal thoughts

## Reformulating

Dr: Paroxetine

P: **Fluoxetine**

Dr: Ah Fluoxetine

Dr: Right oh that's right so it's that it's gone back up to 130

P: **150**

Dr: 150

Dr: Who's your key worker there do you know

P: **Err the person who comes to see me?**

Dr: Yeah the person you see most often I suppose

## Extending

Dr: Yeah well as um shall we just um re re-start where we were we just commencing starting the interview when we um coz we see you was it couple of months three months

P: **Since I saw you?**

Dr: Yeah when was the last time I saw you

Dr: Can you remember what you are on five

P: No

Dr: Or

P: **10 milligrams?**

Dr: 10 milligrams

## Combinations

Dr: Have you experienced this sensation in the past

P: **Have I what?**

Dr: If you have experienced this sensation in the past

Dr: So how are the headaches have they changed at all

P: **What do you mean changed** I got a headache now

Dr: Have they got worse or are they getting

Dr: Are you suspicious are you suspicious of people

P: **Suspicious?**

Dr: Paranoid

P: **Jealous?**

Dr: Jealous yeah

## Combinations

Dr: Aaa so have you had any more thoughts about studying

P: **What music?**

Dr: Well you just you need to come up with a few ideas about what you might study

Dr: Do you do you really feel it or is it a sensation

P: **Is it what I'm thinking is that what you mean?**

Dr: No is it just err the mind playing tricks on you or is it something

Dr: That's right I don't think we've actually booked another time for the Clozapine clinic

P: **Have we already?**

Dr: No I don't think we have



# Method

- Define features manually, then extract automatically:
  - Linguistically/observationally informed:
    - wh-question words, closed-class repair words
    - repeated fragments
    - ...
  - Brute force:
    - *all* the words used (unigrams)
    - patient only, to avoid doctor-specificity
- Train machine learning classifiers to detect NTRIs/P2Rs
  - Supervised classification (SVMs)
- 138 dialogues, c.44,000 speaker turns (c.21,000 by patient)
  - 567 NTRIs (159 patient), 830 P2Rs (262 patient)
  - 5-fold cross-validation

# Features

- Full feature set (one row, one proportional):

Feature	Description
Speaker	Doctor, Patient, Other
NumWords	Number of words in turn
OpenClassRepair	Contains <i>pardon, huh</i> etc
WhWords	Number of wh-words (e.g. <i>what, who, when</i> )
Backchannel	Number of backchannels (e.g. <i>uh-huh, yeah</i> )
FillerWords	Number of fillers (e.g. <i>er, um</i> )
RepeatedWords	Number of words repeated from preceding turn
MarkedPauses	Number of pauses transcribed
OverlapAny	Number of portions of overlapping talk
OverlapAll	Entirely overlapping another turn

## Results - balanced data

- Repair detection, balanced (i.e. small!) dataset:

Target	Features	Accuracy (%)
NTRI	Repeated proportion	65.9
NTRI	All high-level	78.1
NTRI	All unigrams	78.4
NTRI	All features	80.4
P2R	Repeated proportion	64.5
P2R	All high-level	75.7
P2R	All unigrams	77.2
P2R	All features	79.9

## Results - balanced data

- Repair detection, balanced dataset, patient only:

Target	Features	Accuracy (%)
NTRI	Repeated proportion	61.2
NTRI	All high-level	83.4
NTRI	All unigrams	82.4
NTRI	All features	86.3
P2R	Repeated proportion	61.5
P2R	All high-level	78.5
P2R	All unigrams	77.1
P2R	All features	79.8

## Results - balanced data

- Repair detection, balanced dataset, patient only:

Target	Features	Accuracy (%)
NTRI	Repeated proportion	61.2
NTRI	All high-level	83.4
NTRI	All unigrams	82.4
NTRI	All features	86.3
P2R	Repeated proportion	61.5
P2R	All high-level	78.5
P2R	All unigrams	77.1
P2R	All features	79.8

- But this is a sparse phenomenon (0.8% of turns)

## Results - raw data

- Repair detection, raw (unbalanced) dataset:

Target	Features	F (%)	P (%)	R (%)
NTRI	High-level	27.3	36.0	22.3
NTRI	All	32.9	38.9	29.2
P2R	High-level	24.2	32.7	19.3
P2R	All	30.9	37.5	26.5

- Repair detection, raw (unbalanced) dataset, patient only:

Target	Features	F (%)	P (%)	R (%)
NTRI	OCRProportion	35.8	85.7	22.6
NTRI	High-level	41.4	42.8	40.6
NTRI	All	44.0	44.9	43.6
P2R	OCRProportion	19.6	56.4	11.8
P2R	High-level	31.6	36.2	28.4
P2R	All	35.4	43.8	30.3

# Next steps

- We're ignoring non-transcript features
  - Intonation
  - Non-verbal behaviour
- We're ignoring dialogue *context*
  - Human annotators rely on subsequent turn
  - Presence of P3OISR
  - Some similar features to NTRIs (repetition etc)
- Joint problem:
  - Some similarity with decision detection
  - Fernandez et al (2007); Bui & Peters (2010)

# Outline

- 1 Background
  - Doctor-patient communication and adherence
  - Repair
- 2 Repair in patient-doctor communication
  - Method
  - Results
- 3 Automatic detection of repair
  - Approach & previous work
  - Results
  - Next steps
- 4 Automatic prediction of adherence
  - Approach & previous work
  - Results
  - Next steps



## Method

- Apply the same approach to classifying entire dialogues
  - and therefore individual patients
- 125-138 dialogues only!
  - 37 associated with low subsequent adherence
  - 5-fold cross-validation
- Features normalised per dialogue, per word, per turn
- Lexical unigram feature space is very large ...
  - use correlation to find most predictive
  - patient only, to avoid doctor-specificity

# Features

- Feature set used (one each for Doctor, Patient, Other):

Feature	Description
Turns	Total number of turns
Words	Total number of words spoken
Proportion	Proportion of talk in words
WordsPerTurn	Average length of turn in words
WhPerWord	Proportion of wh-words
OCRPerWord	Proportion of open class repair initiators
BackchannelPerWord	Proportion of backchannels
RepeatPerWord	Proportion of words repeated
OverlapAny	Proportion of overlapping talk
OverlapAll	Proportion entirely overlapping other turn
QMark	Proportion containing question intonation
TimedPause	Pause of more than c.200ms (where marked)

## Results - raw data

- Adherence prediction, raw (unbalanced) dataset:

Features	F (%)	P (%)	R (%)
Baseline (all)	44.8	28.9	100
High-level	40.4	28.3	78.6
+ repair features	40.4	28.3	78.6
All features	45.4	29.6	100
Best features ( $\geq 10\%$ )	71.1	62.9	88.9
Best features ( $\geq 3$ )	86.2	89.4	84.8

## Results - raw data

- Adherence prediction, raw (unbalanced) dataset:

Features	F (%)	P (%)	R (%)
Baseline (all)	44.8	28.9	100
High-level	40.4	28.3	78.6
+ repair features	40.4	28.3	78.6
All features	45.4	29.6	100
Best features ( $\geq 10\%$ )	71.1	62.9	88.9
Best features ( $\geq 3$ )	86.2	89.4	84.8

- “Best” is surely unlikely to generalise
- Need to find a more general representation of content

## Topics - predicting non-adherence

- Words chosen reflect some topical content:

air	fill	mates	simply
anyone	finished	monthly	sodium
balanced	fish	mouse	stable
bleach	flashbacks	nowhere	stock
build	grass	pains	symptoms
building	grave	possibly	talks
busy	guitar	pr	teach
challenge	h	recent	terminology
chemical	hahaha	removed	throat
complaining	lager	ri	virtually
cup	laying	schizophrenic	was
dates	lifting	sensation	wave
en	lucky	sickness	worse

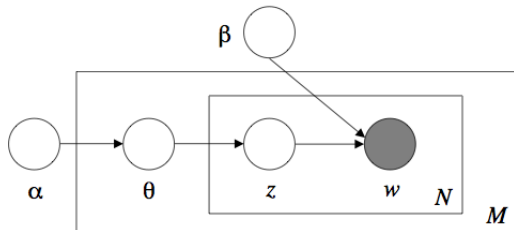
## Topics - predicting PEQ overall

- Different content with patient evaluation:

20th	electric	onto	sometime
ages	energy	overweight	son
angry	environment	oxygen	standing
anxiety	experiencing	packed	stomach
background	facilities	percent	suddenly
bladder	friendly	personally	sundays
booked	helps	picture	suppose
boy	ignore	played	table
broken	immediately	programs	team
bus	increased	progress	television
certificate	irritated	provide	thursdays
dead	kick	public	troubles
deep	later	quid	uhhm
drunk	lee	radio	upsetting
earn	loose	realised	walks
eeerrrr	low	reply	watchers

## Using topics as features

- Existing manual definition of 20 “topics”
  - Medication, side-effects, treatment, management
  - Symptoms, health
  - Daily activities, living situation, relationships, ...
  - (but not annotated over sufficient data)
- Can we learn topics from the data?
- Latent Dirichlet Allocation (Blei et al, 2003)



# LDA topics

- Most probable words for learned topics:

Topic 0	feel low alright mood long drug feeling tired time confidence coming
Topic 4	voices pills mood cannabis telly voice shaking chris control inside ma
Topic 5	letter health advice letters council copy send dla cpn problems housin
Topic 7	church voice voices hear medication sister bad hearing taking felt nev
Topic 9	school children kids back september oclock gonna phone social son w
Topic 10	weight months medication stone risk lose eat write gp hasnt exercise
Topic 11	place support work centre gotta job stress feel psychologist theyll cor
Topic 12	door house police thought ring knew worse wall hadnt sat coming fea
Topic 13	doctor alright years nice ill anxious write long sit eye heart ring lovely
Topic 14	drug taking milligrams hundred doctor night time medication voices
Topic 15	sort medication work drugs kind team issues drink alcohol things sup
Topic 16	mum place brother tablets died dad depot house meet money lives da
Topic 17	people life drug make care lot friends dry camera live cope thing can
Topic 18	alright house drink drinking money alcohol god drugs living basically
Topic 19	kind day time remember side weeks blood hospital appointment case



## Results - with topic features

- Include topic weight per dialogue as features
- Adherence prediction, raw (unbalanced) dataset:

Features	F (%)	P (%)	R (%)
Baseline (all)	44.8	28.9	100
High-level	40.4	28.3	78.6
+ repair features	40.4	28.3	78.6
All features	45.4	29.6	100
+LDA features	56.3	47.3	69.8
Best features ( $\geq 10\%$ )	71.1	62.9	88.9
Best features ( $\geq 3$ )	86.2	89.4	84.8

## Results - with topic features

- Include topic weight per dialogue as features
- Adherence prediction, raw (unbalanced) dataset:

Features	F (%)	P (%)	R (%)
Baseline (all)	44.8	28.9	100
High-level	40.4	28.3	78.6
+ repair features	40.4	28.3	78.6
All features	45.4	29.6	100
+LDA features	56.3	47.3	69.8
Best features ( $\geq 10\%$ )	71.1	62.9	88.9
Best features ( $\geq 3$ )	86.2	89.4	84.8

- Perhaps a more generalisable result?

## Next steps

- Manually annotate for topic
- Investigate discourse-based and (semi-)supervised LDA
- Investigate role of other dialogue phenomena
  - Doctor-led repair
  - Non-verbal communication