

School of Electronic
Engineering and
Computer Science

MSc. Big Data Science
Thesis Project Report
2017

**Interactional and
Linguistic Analysis
for Computationally
Diagnosing
Alzheimer's Disease**

Claire Mary Kelleher
Supervisor: Dr. Matthew Purver

Disclaimer

This report, with any accompanying documentation and/or implementation, is submitted as part requirement for the degree of MSc. in Big Data Science at Queen Mary University of London. It is the product of my own labour except where indicated in the text. The report may be freely copied and distributed provided the source is acknowledged.

Abstract

Background: On top of memory loss and linguistic impairment, changes in behaviour and decreased interactional skills in conversation are also symptoms for Alzheimer's Disease (AD). Relatively few studies have used computational techniques to investigate the diagnosis predicting power of interactional symptoms. Automatic diagnosis of AD aids assessment and allows for earlier diagnosis. Interactional features are linguistically and culturally independent, allowing the automation to be applied across languages and borders.

Overview: This project investigates what interactional features between an invigilator and patient can be utilized to predict AD and how they perform against and in addition with the proven top predicting linguistic features (non-interactional). This was done by re-encoding the optimal predicting non-interactional features (Non-IF's) and the classifier used to find them (based on [Fraser et al. 2016]). This resulted in similar classifier performance scores as [Fraser et al. 2016] allowing this result to be used as an evaluation baseline to compare how the classifier performs when interactional features were added. Interactional features (IF's) that are known symptoms of AD such as turn-taking, filler term frequency and trailing-off mid-sentence, were encoded. Non-IF's & IF's were compared by ranking the features based on classifying weight. IF's were hierarchically grouped by 3 umbrella features; 'Fillers', 'Unintentional Silence' and 'Dialogue'. Correlation analysis was carried out on the interactional features as a sense check and to investigate the direction in which each variable correlates with the diagnosis.

Evaluation Method: The predictive weight scores (ANOVA F-Value) were calculated for the top LF's. For hypothesis purposes, it was assumed that these were the top predicting variables. The null hypothesis was there is no change in these top predicting variables. The alternative hypothesis was that the top predicting variables would include both IF's and LF's. Where, because IF's with a higher weight have replaced lower weighted LF's in the rank, also leads to an improvement in classifier performance.

Results:

- The final top predicting features were made up of 52.2% interactional and 47.8% linguistic with classifier precision improving by 4.96%. (See Appendix Table 3)
- Dialogue based features, alongside lack of clarity between the examiner and participant, resulted in higher predictive weight almost 3 times stronger than the other variables.

Conclusion: Since IF's replaced over half of the originally all LF's top predicting features regarding diagnosis predicting weight, as well as causing an improvement in the classifier's accuracy, one can conclude that encoding interactional features, in particular dialogue based features that represent the amount of involvement the invigilator is required to manifest throughout the session, in addition to non-interactional features, can assist in computationally classifying AD.

Acknowledgments

A special thank you to my supervisor and lecturer, Dr. Matthew Purver, for not only introducing me to the fascinating subject of NLP, but for giving me this opportunity to explore the area and expand my understanding of it. Your guidance was very much appreciated throughout the year.

Thank you to all my girls, for not distracting me too much. My housemates in particular, for being so understanding when I occupied the entire kitchen to study.

And of course, thank you Mam and Dad, for your love and unconditional patience.

Table of Contents

1. Background.....	7
1.1 Dementia & Alzheimer’s Disease.....	7
1.2 Project Introduction.....	8
1.3 Automatic Classification of Alzheimer’s Disease.....	9
2. Literature Review & Motivation	11
2.1 Literature Review.....	11
2.2 Motivation for Investigating Interactional Features	15
2.3 Hypothesis.....	16
3. Materials & Methods	17
3.1 Dataset & Target.....	17
3.2 Method.....	20
4. Implementation	22
4.1 TalkBank Manual: CHAT & CLAN	22
4.2 Python	22
4.3 Pre-Processing	22
4.4 Feature Extraction: Non-Interactional Features	23
4.4.1 Factor 1. Semantic Impairment	24
4.4.2 Factor 2. Syntactic Impairment.....	27
4.4.3 Factor 3. Information Content Impairment.....	27
4.5 Feature Extraction: Interactional Features.....	29
4.5.1 Umbrella Feature 1. Fillers	30
4.5.2 Umbrella Feature 2. Unintentional Silence	32
4.5.3 Umbrella Feature 3. Conversation.....	33
4.6 Non-Interactional Features Selection.....	35

4.7 Classifier Build	37
4.8 Interactional Features Selection	39
4.9 Diagnosis Correlation Analysis of Interactional Features	40
5.Evaluation Methods	43
6.Results.....	45
6.1 Discussion	46
6.2 Conclusion.....	47
6.3 Further Work.....	48
7. Bibliography	50
8. Appendix	54

1. Background

1.1 Dementia & Alzheimer's Disease

The word 'dementia' describes a set of symptoms that may include memory loss and difficulties with thinking, problem-solving or language. Dementia is caused when the brain is damaged by diseases, such as Alzheimer's disease (AD) or a series of strokes. AD is the most common cause of dementia, but not the only one. The specific symptoms that someone with dementia experiences will depend on the parts of the brain that are damaged and the disease that is causing the dementia [Alzheimer's Society - What is Dementia? 2017]. Dementia is a terminal condition and affects millions of people worldwide. Symptoms of dementia include memory loss, confusion and problems with speech and understanding. Due to Dementia having a high prevalence and level of associated morbidity, it has become an urgent health and economic issue for the developed world, and a rapidly growing threat in developing countries. There are currently 850,000 people with dementia in the UK, with numbers set to rise to over 1 million by 2025. Due to the general rise of elderly individuals in the population, this will soar to 2 million by 2051. [Alzheimer's Society - Facts for the media 2017]

AD is a degenerative brain disorder and the most common type of dementia, affecting 62% of those diagnosed. Other types of dementia include; vascular dementia, affecting 17% of those diagnosed, mixed dementia affecting 10% [Alzheimer's Society - Facts for the media 2017]. The risk of developing AD becomes greater with age. AD can refer to a set of symptoms in different cognitive and linguistic domains, and characteristically, these symptoms are persistent and progressive, causing a deterioration of skills and knowledge. The domains affected are memory, executive functions, language, visual-spatial processing, personality, general behaviour and interactional skills [Svennevig and Lind 2016]. This

project focuses on exploring the effects that AD has on both the impairment of linguistic and interactional skills using computational techniques.

1.2 Project Introduction

After memory loss, language impairment and changes in a patient's behaviour are part of the top symptoms of AD. A number of studies have already been carried out using computational techniques on linguistic features in order to identify AD [Fraser et al. 2016, Croisile et al. 1996, Guinn and Habash 2012, Bucks et al. 2000]. For example, [Fraser et al. 2016] demonstrated high end accuracy automatically identifying AD from short narrative samples elicited with a picture description task and [Croisile et al. 1996] provides comparative information between AD patients and healthy elderly subjects about lexical, syntactic, and semantic aspects of oral and written picture descriptions. However relatively few studies have been carried out where the features that are encoded represents the patients interactional behaviour. It is hypothesised that by including this level of symptoms that has not been investigated before in the AD classification process, the performance of the classifier will improve.

The aim of this project was to encode a set of interactional features (IF) that represent common interactional behaviour changes in AD patients such as difficulties following a conversation or finding the right word for something [Alzheimer's Society - What is Dementia? 2017], hesitation, restlessness [Alzheimer's Society - Behaviour Changes 2017] (eg. conversation restarts, filler terms, pauses) and to investigate the predictive power of these features in classifying the different AD severities both on their own and combined with the non-interactional features that are already proven to be useful [Fraser et al. 2016] in making these predictions. An investigation into the correlation between the interactional features and the diagnosis variables was also carried out in order to investigate in which features most correlate with it.

1.3 Automatic Classification of Alzheimer's Disease

Due to the general rise of elderly individuals in the population, the number of AD and dementia patients through the UK and the world is increasing. This leads to a higher demand of care from the healthcare system where staffing levels are not able to keep up becoming an urgent health and economic issue for the developed world. The automatic classification of AD can reduce these pressures by assisting the practitioner's assessment of patients, acting as a second opinion, leading to an earlier diagnosis of the illness which ultimately increases efficiency and reduces costs, benefitting both the patients and the healthcare system.

A number of studies have been carried out to automate the diagnosis of AD and other dementias, using a variety of classification features such as SPECT¹ images [Horn et al. 2009], MR scans [Klöppel et al. 2008] and linguistics elicited from both narrative speech [Fraser et al. 2016] and spontaneous speech [Lopez-de-Ipiña et al. 2015]. Patient's linguistic features are the more cost effective and less invasive to record and analyse than brain imaging (SPECT & MR).

The above features, however, do not measure a change in patients interactional and dialogue behaviour, which is a key symptom in diagnosis AD and other dementia. [Alzheimer's Society - Behaviour Changes 2017]. The encoding of this layer of symptoms is not only estimated to be of similar low invasive and cost levels of linguistic based analysis, but is also less dependent on specific languages and cultures. For example, it has been shown in [Dingemanse et al. 2013] that the clarification or filler term 'Huh?' is a universal word likely to be attested in similar form in all natural spoken languages and so can be said that this interjection can transcend between cultures and languages. If interactional features

¹ Single-photon emission computed tomography

such as the process of examiner-participant clarification and other interjections are encoded to classify AD, it can be said that this classification process can be applied to different countries and for a range of language tasks.

2. Literature Review & Motivation

2.1 Literature Review

In “Linguistic Features Identify Alzheimer's Disease in Narrative Speech” [Fraser et al. 2016], Fraser, Meltzer and Rudzicz demonstrated state-of-the-art accuracy in automatically identifying AD from short narrative samples elicited with a picture description task, and to uncover the salient linguistic factors with a statistical factor analysis. Data was derived from the DementiaBank corpus [DementiaBank 1994], from which 167 patients diagnosed with “possible” or “probable” AD provide 240 narrative samples, and 97 controls provide an additional 233 files from 97 speakers. This corpus consists of narratives that are segmented into utterances and manually annotated with filled pauses, paraphasias, and unintelligible words however, Fraser et al. opted to use the Stanford Parser² and Stanford Tagger³ to annotate the data and used the manually annotated tags from DementiaBank corpus to test the performance of the Stanford Tagger on the data resulting in an accuracy of 85.4% on the control data, and 84.8% on the AD data (over the entire data sample).

Using machine learning feature selection methods, Fraser et al. have carried out a significant exploration into what features, the optimal number of features and what combination of features best classify the different AD. The maximum highest accuracy score (81%) for their model (logistic regression) was received when 35 features were inputted. This was followed up with an exploratory factor analysis where factor were found using a promax oblique rotation. This analysis was carried out on the top 50 of the 370 features. This number of variables were chosen because it was found that there was not much change

²Version 2010-11-30 <https://nlp.stanford.edu/software/lex-parser.shtml>

³Version 2015-01-29 <https://nlp.stanford.edu/software/tagger.shtml>

in the classifiers accuracy score until more than 50 features were added, where there was then a sharp drop in accuracy score. The factor analysis resulted in finding four factors or groupings of features: semantic impairment, acoustic abnormality, syntactic impairment, information impairment.

Fraser et al. obtained classification accuracies of over 81% in distinguishing individuals with AD from those without based on short samples of their language on the Boston 'Cookie-Theft' picture description task [Goodglass and Kaplan 1983]. Using factor analysis, four clear factors emerged from the study: semantic impairment, acoustic abnormality, syntactic impairment, and information impairment. A large number (370) of features were considered to capture a wide range of linguistic phenomena. The papers classifier accuracy and build was also used as a baseline for the evaluation when comparing both the interactional and non-interactional features.

The same picture task is used in "A Comparative Study of Oral and Written Picture Description in Patients with Alzheimer's Disease" [Croisile et al. 1996] as the one used in [Fraser et al. 2016]. The participants in this study consisted of 22 patients with AD and 24 healthy elderly subjects. The purpose of the study was to provide comparative information about lexical, syntactic, and semantic aspects of oral and written picture descriptions in AD patients and healthy elderly subjects. Croisile et al. analysed the similarities and differences of oral and written descriptions by comparison of the results obtained in each group and identified specific impaired features of description processing in AD patients by making an intergroup comparison of the results obtained for each task. The result was that AD patients had a significant reduction of all word categories, which, similarly to controls, was more pronounced in written than in oral texts and in sum, AD descriptions were always shorter and less informative than control texts. [Fraser et al. 2016] based their studies definition of an information unit on [Croisile et al. 1996]. An information unit is a means to measure the information content that the participant described based on a list created by [Croisile et al.

1996]. The list consisted of 23 information units in four key categories: subjects, places, objects, and actions. For example, the three subjects were: the boy, the girl, and the woman. If a participant mentioned mother or female, this would count as a mark for the woman information unit.

Some other research investigates how interactional features are effected by dementia in Norwegian speaker who are multilingual such as “Dementia, interaction, and bilingualism: An exploratory case study” [Svennevig and Lind 2016]. The study presents an exploratory, clinical linguistic case study of one bilingual speaker diagnosed with probable dementia of the Alzheimer type in two conversational contexts, English and Norwegian. The study explores his speech production in the two languages, focusing on one case where the participant displays problems of achieving progressivity of talk and his methods in which he searches for ways of continuing his turn at talk. The study investigates turn-taking and take the most crucial word of an utterance to be the lexical verb as so much of the semantic and syntactic structure of the utterance, hence also the interpretation of the utterance, depends on the choice of the lexical verb. Without lexical verbs, utterance interpretation is very challenging, even in context. The study looks at the participant's word-finding difficulties manifested as a lack of progressivity of the talk and investigates the way in which the participant solves these difficulties. For example, solutions for this difficulty would include the speaker saying semantically meaningless nouns (e.g. “thing”) as a substitute for the intended word or fillers (e.g. “uh”, “em”). Sometimes, the act of searching is not solely caused by the speaker but as a joint effort between the different conversation participants. The speaker may include an invitation by the other party in the search, for instance by gazing at him or her, or by explicitly appealing for assistance. The data collected included responses to formal cognitive and linguistic tests, as well as responses to a questionnaire on functional communication and recordings of more or less spontaneously produced speech (elicited narratives and conversation). The study focuses primarily on the conversational data, while

using some of the other data as background for the description of the participant. Although [Svennevig and Lind 2016] is an exploratory study focusing on multi-linguists in Norway, the interactional features mentioned above, plus a number more used in the paper can be said to measure anomia, a symptom of AD. This is the reason behind a number of features were created based around the turn-taking and tokens that may represent interruptions in an utterance and difficulties in progressivity of talk in the hope that they would be strong predictors of AD.

In [Silvast, M. 1991] Aphasia therapy dialogues, Silvast, M. Investigates the interaction between aphasiac patients and the speech therapist in order to investigate the role of the therapist when using conversation as a method for rehabilitating aphasic patients. Aphasia, is the inability to understand or produce speech and has been a proven symptom of AD [Fergadiotis et al. 2011]. These interactions between the patient and the therapist were explored by video-taping a fraction from a therapy session. Six aphasic-therapist pairs served as subjects in the study. A middle five-minute segment of each conversation was extracted for analysis, which focused on the use of interactional space and different communicative functions in therapy conversations. The results showed that during the conversation, therapists had a regulatory role which was manifested in their frequent use of requests for information and clarification. The reason the aphasiacs had a longer speech time is explained by their frequent trouble-indicating behaviour during speaking such as pauses, fillers and repeats.

An investigation into the interjection and filler term ‘Huh?’ was carried by [Dingemanse et al. 2013] in “Is ‘Huh?’ a Universal Word? Conversational Infrastructure and the Convergent Evolution of Linguistic Items” to determine if it is firstly, universal and secondly, a word. In support of the first investigation, it was shown that the similarities in form and function of the interjection across languages are much greater than expected by chance.

The second investigation of the study shows that 'Huh?' is a lexical, conventionalised for that has to be learnt, unlike grunts or emotional cries. Both investigations were carried out across ten different languages with 196 instances of the interjection for other-initiated repair (OIR) were collected through video recordings.

It was concluded that 'Huh?' is a universal word, not because it is innate but because it is shaped by selective pressures in an interactional environment that is shared by all languages. This environment is the other-initiated self-repair which then leads to by [Dingemanse et al. 2013] suggesting that conversational infrastructure can drive the convergent drive the convergent cultural evolution of linguistic terms.

2.2 Motivation for Investigating Interactional Features

There has been a number of studies carried out using computational techniques on the linguistic based symptoms of dementia [Fraser et al. 2016][Croisile et al. 1996] however there has been relatively little research done on symptoms that fall under what can be classified as interactional symptoms such how a patient solve anomia (word retrieval problems) and what solution path they take. As mentioned in the literature review, individuals with dementia of the Alzheimer's type often experience anomia [Svennevig and Lind 2016] leading to the use of filler terms. Hesitant speech also increases with the severity level of the dementia [Szatloczki et al. 2015], again a known situation where the speaker either intentionally or unintentionally utters filler terms or pauses, leading to a lull in the conversation flow, effecting the interaction between the conversation participants. Other non-linguistic, interactional features that become more prominent as AD progresses include the decline in the patient's ability to concentrate, making it difficult for the patient to complete tasks and to follow conversation. [Alzheimer's Society - What is Dementia? 2017] This can phenomenon can be encoded by analysing the length of an AD patient's answers or the if

they needed to clarify the task at hand by asking the examiner questions, which has been proven to happen with aphasic patients [Silvast, M. 1991].

By investigating these interactional features, another layer of dementia symptoms can be added to the analysis which have not been looked at using computational techniques and NLP methods before. It is theorised that by adding in more variables that represent another layer of dementia symptoms, the diagnoses of dementia types using computational techniques should have a higher accuracy result.

[Dingemanse et al. 2013] claims that the interjection 'Huh?', used for clarification purposes, is shaped by selective pressures in an interactional environment. This project explores the use of other interjections such as pauses that may have been manifested from the interactions between the patient and examiner and aims to use them to represent the interactions between both parties.

2.3 Hypothesis

It is hypothesised that the addition of the interactional features such as the ones investigated in [Silvast, M. 1991] and [Svennevig and Lind 2016] will improve the classification of AD than using non-interactional features only, such as those in [Fraser et al. 2016]. This is because the additional interactional features represent another level of symptoms used to diagnose AD as opposed to only taking into account the patient's linguistic symptoms.

3. Materials & Methods

3.1 Dataset & Target

Following the studies carried out by [Fraser et al. 2016], the data used was derived from the DementiaBank (Pittsburgh) corpus [DementiaBank 1994], downloaded in June 2017, which is part of the larger TalkBank project [TalkBank 2007]. The Alzheimer Research Program at the University of Pittsburgh [Becker et al. 1994] collected the data between 1983 and 1988. Participants were referred directly from the Benedum Geriatric Centre at the University of Pittsburgh Medical Centre, and others were recruited through the Allegheny County Medical Society, local neurologists and psychiatrists, and public service messages on local media.

The Pittsburgh corpus that was used for this project consisted of 550 transcripts. Each transcript is a recording between one participant and the examiner where the examiner used the “Cookie Theft” picture description task from the Boston Diagnostic Aphasia Examination [Goodglass and Kaplan 1983] on the participant. This task instructs the examiner to show the picture to the patient and say, “Tell me everything you see going on in this picture.” This test consists of the examiner telling the participant to describe everything that they see in the picture that is being shown to them. The picture is a cartoon image of a scene set in a kitchen, where a boy is robbing cookies while his mother is cleaning the dishes the examiner is permitted to encourage the patient to keep going if they do not produce very many words. This leads to the interactions and pauses in conversation between the participant and the investigator that can be investigated for this project. Verbal picture description is one of the most sensitive tests for detecting language disorders in early AD [Carozza 2015] and is the

reason this task is used in a number of studies to investigate AD and other dementias. [Croisile et al. 1996, Giles et al. 1996, Bird 2000]

Each file from the Pittsburgh corpus is one transcription of a session between an examiner and one participant, and so represents one distinct patient with one diagnosis (target variable). Variables for this analysis were created based on one value (average or otherwise) that represents the feature in that particular transcript. In other words, the data that the model was trained and tested on was at transcript level and each feature was represented by a single number for each file. As mentioned earlier, there was one target value per transcript so this was already on the correct level for analysis. The five diagnoses (target) labels were 'MCI' (Mild Cognitive Impairment), 'Memory', 'Possible AD', 'Probable AD', 'Vascular' (Vascular Dementia), where these are approximately in order of severity respectively.

MCI is an intermediate stage between the expected cognitive decline of normal aging and the more-serious decline of dementia. It can involve problems with memory, language, thinking and judgment that are greater than normal age-related changes. [MFMER 2016] There were 43 participants in the study with MCI. 'Memory' classifies patients that are experiencing memory loss only and at a rate only slightly greater than that of the changes that would be expected for a person at that age. The data only consisted of 3 'Memory' diagnosis. Possible AD and Probable AD are the different severities of Alzheimer's disease with 'Probable' being more severe than 'Possible'. There were a total of 237 'Probable AD' transcripts and there was only 21 transcripts that recorded patients that were classified with the 'Probable AD' diagnosis. Vascular dementia is caused when blood flow to the brain is reduced. It is the second most common cause of dementia, after AD [Alzheimer's Society - Vascular Dementia 2017]. There were only 5 patients in the DementiaBank data with Vascular dementia. Due to the dataset being unbalanced, precautions were taken to ensure results could be trusted. This included the use of cross-validation when training the model

as well as taking into account the recall score (and F1 score) when looking at classifier accuracy scores. Different cuts of the data (e.g. 'Control' against 'Probably AD' only) and different groupings of the label's (e.g. group any non-AD label as 'Other') were also investigated.

Each file (one participant-examiner medical session) is speech sample was recorded then manually transcribed at the word level following the TalkBank CHAT (Codes for the Human Analysis of Transcripts) protocol [MacWhinney, B. 2000], which is discussed in more detail in the follow section. Narratives were segmented into utterances and annotated with filled pauses, paraphasias, and unintelligible words. CHAT is the standard transcription system for the TalkBank and CHILDES (Child Language Data Exchange System) Projects. All of the transcripts in the TalkBank databases are in CHAT format. These annotations were used to encode and analyse the interactional features of the patients in this study where as they were not used for the non-interactional (linguistic). Instead, the transcripts were passed through the Stanford Tagger⁴ and Stanford Parser⁵. This is because as mention in the literature review, [Fraser et al. 2016], did so and by running the transcripts through the same parser (although a more recent version of the Stanford tagger and parser were used for this project) as [Fraser et al. 2016], this project was allowed compare the results for the non-interactional features to those of [Fraser et al. 2016]. More information on other pre-processing steps that had to be taken are discussed in the following sections.

3.2 Method

The following experiments first attempts to replicate the non-interactional features used in the classification experiments from [Fraser et al. 2016]. A new set of interactional features

⁴Version 2017-06-09 <https://nlp.stanford.edu/software/tagger.shtml>

⁵Version 2017-06-09 <https://nlp.stanford.edu/software/lex-parser.shtml>

(IF's) were introduced and compared against the [Fraser et al. 2016] variables regarding AD classification accuracies. These experiments are carried out in order to investigate how a change in person's interactional features in conversation can be used to better classify the different diagnoses of AD using NLP methods and other computational techniques. This means that the main requirement for this project was to produce a set of encoded variables that represent both linguistic and interactional phenomena. Once encoded, these variable sets can be compared against each other regarding their classifying weight towards predicting the target variable (diagnosis). The target variable was the five classes mentioned previously of the dementia diagnosis as well as the 'Control' group. The five diagnoses labels are 'MCI' (Mild Cognitive Impairment), 'Memory', 'Possible AD', 'Probable AD', 'Vascular' (Vascular Dementia), where these approximately in order of severity respectively.

Using machine learning feature selection methods and the same data from DementiaBank, [Fraser et al. 2016] have found the optimal number (35) and strongest AD predicting features that were factored into are semantic, syntactic, information content and acoustic impairment based. For this project, it was decided to re-build the features from the semantic, syntactic and information impairment factors (in other words, ignoring the acoustic factor) from [Fraser et al. 2016] and use them to define the non-interactional features (Non-IF's). A total of 29 Non-IF's were encoded. Details on each feature and the reasoning behind the algorithm used to encode each one can be found in the following section. A number of classifiers trained on the DementiaBank data using the Non-IF's as variables and the, including Decision Tree, Logistic Regression, K-Nearest Neighbour and Naive Bayes. Different labelling's of the target variable were also used. For example, one cut of the data counted any non-AD targets as 'Other'.

This project encoded Non-IF's and trained and tested a classifier to replicate that of [Fraser et al. 2016]. This outputted a similar classifier accuracy scores to that of [Fraser et

al. 2016] and so allowing this score to be used as an evaluation baseline to compare how the classifier performs when other features were inputted.

A number of interactional features that are known symptoms of AD such as pauses and filler term frequency [Svennevig and Lind 2016] and trailing-off mid-sentence [Alzheimer's Society - What is Dementia? 2017] and different aspects of dialogue such as turn-taking ratio [Silvast, M. 1991], were encoded. Different aspects of these umbrella IF's were chosen in a more investigative approach such as investigating if the location of the filler played a role in predicting AD or the type of word (POS tag) that patients tending to forget correlated with the diagnosis.

The two groupings of features (non-interactional and interactional) were then compared by firstly, ranking the features based on classifying weight (ANOVA F-Value) and secondly, running the top predicting features from this set through the classifier built earlier and comparing the change in accuracy. Correlation analysis was carried out on the interactional features as a sense check and to investigate the direction in which each variable correlates with the diagnosis.

4. Implementation

4.1 TalkBank Manual: CHAT & CLAN

The first step of implementation was to investigate the TalkBank transcripts and have an understanding of what the data represented. In order to do this, an understanding of the TalkBank [TalkBank 2007] manual was required. There are three parts to the overall TalkBank manual. Part 1 describes the CHAT transcription system, part 2 describes the CLAN (Computerized Language Analysis) analysis programs and part 3 describes the segments of the CLAN program that perform automatic morphosyntactic analysis.

4.2 Python

Python (3.6) programming language was used to encode the variables. Python has a range of libraries that proved incredibly useful for this project. The libraries that were utilised included 'pandas', 'numpy', 'sklearn', 'scipy' and others. NLTK, Python's platform for working with human language data was also utilised.

4.3 Pre-Processing

The pre-processing step for encoding the Non-IF's differed from that of the IF's. This is because when [Fraser et al. 2016] encoded their features, they only kept the word level transcription and the utterance segmentation. In other words, they discarded the morphological analysis, disfluency annotations, and other associated information that TalkBank had annotated. Whereas, variables created to represent the IF's needed these annotations as they represent interactional phenomena, both of the participant and between the participant and the invigilator.

For the Non-IF extraction, even after removing there CHAT morphological annotations, there were a still a number of unneeded symbols within the tokens of the transcript. Since the participants words were only needed for the analysis, these unwanted symbols and digits were removed using regular expression. The remaining tokens were lemmatized and passed through the Stanford Tagger and Stanford Parser (versions stated in earlier 'Method' section).

4.4 Feature Extraction: Non-Interactional Features

As stated earlier, [Fraser et al. 2016] have found the optimal number (35) and strongest predicting features that are linguistic and information content based (non-interactional), that best classify AD. Factor analysis was on the top 50 features (out of 370) to finding underlying groupings of the features. The factor analysis resulted in finding four factors or groupings of features: semantic impairment, acoustic abnormality, syntactic impairment, information impairment. The encoding of the Non-Interactional Feature (Non-IF's) was based on the features within the non-acoustic factors. I.e. semantic impairment, syntactic impairment, information impairment. [Fraser et al. 2016] found these factors were found by carrying out a promax oblique rotation. For this project, the features with the most significant factor loading scores (greater than 0.3) were chosen to be encoded to represent the Non-IF's (Appendix Table 7). This was so the factor interpretations (groupings) could be carried forward and used through out the analysis.

The following are the features encoded for this project to represents the Non-IF's grouped by their relative factors. A brief explanation on how each variable was encoded, including, if needed, how they were normalised and the name given to the variable for the analysis in brackets. It was not the aim of this project to re-build the features from [Fraser et al. 2016] exactly as they did but to use them more as a strong guidance with the aim of building a classifier that resulted in similar accuracies to that of [Fraser et al. 2016] in order

to compare the change in accuracy when the IF's were included. Regarding the POS (Stanford) tags, a dictionary was created for each script that counted the frequency of each POS tag. Each POS count was normalised by the total number of words spoken by the participant.

For the non-interactional features, all variables were encoded at participant level, and the below features of the invigilator were not taken into account for the Non-IF's extraction. So when the following features state "total tokens in transcript", it refers to total tokens in transcript spoken by the participant.

4.4.1 Factor 1. Semantic Impairment

Pronoun: Noun Ratio: Count of Pronoun / Count of nouns. ('ProNoun_Noun_Ratio')

Adverbs: Frequency of adverbs / total tokens in transcript. ('Adverb_Ct_Total')

Verb frequency: Frequency of verbs / total tokens in transcript. ('Verb_Ct_Total')

Nouns: Frequency of nouns / total tokens in transcript. ('Noun_Ct_Total')

Word length: Script letter count / script word count. I.e. The average word length per script was used to encode this variable. ('Avg_Word_Len')

Honore's Statistic: Honore's Statistic attempts a deeper analysis by accounting for words that are only used once; indicating a higher lexical richness therefore a negative Honore's Statistic suggests low lexical diversity, a known symptom of aphasia and AD. [Fergadiotis et al. 2011] and the exact formula is calculated based on the formula $R = 100 \log N / (1 - V1/V)$ where, R is the Honore's Statistic, N is the total number of tokens, V1 is the words spoken only once and V is the number of lexical words [Honoré 1979]. The lexical words were underlined in each of the samples based on Williamson (2009). Honore's Statistic was represented in this project by creating a dictionary that uses a the list of distinct words in the transcript as its keys and a frequency count in the transcript of each word. A

count of any words that had a frequency of one (i.e. Words used once) per transcript was used to represent the Honore's statistic. ('wordsOnceOverTotal')

Inflected verbs: Inflection is the name for the extra letter or letters added to nouns, verbs and adjectives in their different grammatical forms. Verbs are inflected in the various tenses (-ing,-s, and-ed) [Bloom et al. 1980]. Therefore, to encode this feature, the different verb inflections were counted. A feature for each inflected verb frequency was created. The five encoded verb inflection variables were past tense, gerund (a verb form which functions as a noun), present participle, non-3rd person singular present, 3rd person singular present. As expected there were a lot of correlation between these features so some are removed during the feature selection process. ('Verb_Ct_3sing', 'Verb_Ct_Gerund', 'Verb_Ct_Past', 'Verb_Ct_Past_P', 'Verb_Ct_Past_T')

Average Cosine Distance: This was used to measure participants repetitiveness.[Fraser et al. 2016] measured the cosine distance between each pair of utterances. To encode this variable for this project, each utterance was represented as a term-document vector. TfidfVectorizer, a function that converts a collection of raw documents (in this case utterance) were used to vectorise the utterances. Once the vectorization and transformation was applied to the pair of utterances that the similarity score was being calculated on, a similarity matrix was produced. This was a 2x2 matrix since it was a one utterance being compared against another utterance, where the entries of the matrix were the cosine similarity score. As the entries of diagonal were all equal to 1 (cosine distance), this was able to be used as a sense check because when the utterance was being compared against itself and it obviously going to have cosine similarity score of 1, as both utterance vectors are in the exact same direction (angle between them is 0).

This process was carried out between every participants utterance and their utterance prior to that one (with the exception of first utterance as this did not have an utterance before it) to find the cosine distance between each of them. The sum of these cosine distances was

calculated. This value was divided by the total number of participants utterances in the session in order to find the value for the 'Average cosine distance' variable. ('cos_avg')

Lexical Parsings: Each utterance of the participant was parsed to find its most likely Probability Context Free Grammar (PCFG) tree using the Stanford parser. A variable was created for each of the following parsings was created, which was a count of the occurrence of the parsings throughout the session, divided by the total number of utterances spoken by the participant. It was calculated by traversing through the main parse tree's subtrees. If the root of the subtree was equal to the parent node label (POS tag) in question (before the arrow), then that nodes children nodes were traversed through. If these children nodes labels are equal to the POS tag in questions (after the arrow), a value of 1 was added to the counter for this parsing. This total value for this counter was then divided by the total number of utterance spoken by the participant at the end of the script. The following parsings were encoded and counted per participant per session:

ADVP -> RB ('ADVP_to_RB_norm')

NP -> PRP ('NP_to_PRP_norm')

NP -> DT NN ('NP_to_DT_NN_norm')

4.4.2 Factor 2. Syntactic Impairment

Verbs: Absolute frequency of verbs (base tense) ('Verb_Ct_Base')

Lexical Parsings: As explained earlier, all parsing based features was created by parsing the each of the participants utterances using the Stanford parser and analysing each it's subtree's POS labels and counting all the subtree's within the main tree where the parse in question occurred. In order to normalise, this count was then divided by the total number of utterances spoken by the participant. The same method was used to create these parsing variables as those (ADVP->RB, NP->PRP, NP->DT NN) in the previous semantic factor

section. The following parsing occurrences (one variable each) were counted per participant per session:

ROOT -> FRAG (ROOT_to_FRAG_norm)

VP -> AUX VP ('VP_to_AUX_VP_norm')

VP -> VBG (VP_VBG_norm')

4.4.3 Factor 3. Information Content Impairment

There were two types of features for this factor, keywords and information units. Both related directly to the descriptive content stated by the participant (either they have or have not mentioned certain keywords).

A keyword-based feature is simply a total count of the times a participants mentions a keyword. There were six keywords: 'Window', 'Sink', 'Cookie', 'Curtain', 'Counter' and 'Stool'. These were encoded into six separate features ('Keyword_window', 'Keyword_sink', 'Keyword_cookie', 'Keyword_curtain', 'Keyword_counter' and 'Keyword_stool') in order to investigate if anyone particular keyword had a correlation with the target variables, as well as being summed to create another variable that is the total number of keywords mentioned throughout the session.

An information unit is a measure of how much content the participant has described in the description task. [Fraser et al. 2016] used the following definition of an information unit from [Croisile et al. 1996] which was used to encode this variable for this project. An information unit was a list of words categorized in to four key categories: subjects, places, objects, and actions. The three subjects were: the boy, the girl, and the woman. The two places were the kitchen and the exterior seen through the window. The eleven objects were: cookie, jar, stool, sink, plate, dishcloth, water, window, cupboard, dishes, and curtains. The seven actions or facts were: boy taking or stealing, boy or stool falling, woman drying or washing dishes/plate, water overflowing or spilling, action performed by the girl, woman

unconcerned by the overflowing, woman indifferent to the children. To re-build this feature, utilizing Wordnet's (Python library) 'Synsets', a list of synonyms for each information unit's keyword is created and used to represent an information unit, as opposed to using the exact same list as [Croisile et al. 1996]. If a participant mentions one of the words related to the keywords synonym list, this was a count for the information unit associated with that keyword, as that information unit had been mentioned. The following five keywords were used to build five information units: 'Window', 'Curtain', 'Cookie', 'Sink', 'Girl'. The final variable was a count of how many information units were mentioned. ('Subj_IU_Total').

4.5 Feature Extraction: Interactional Features

For the purpose of this study, the definition of interactional features represents non-linguistic interactions carried out between the participant and the examiner as well as those of the participant and the examiner separately, such as fillers, pauses and conversational interactions such as turn-taking. There were certain lexical aspect to the creation of these interactional features such as analysing the POS tag that appears directly after the occurrence of a filler.

There were a total of 26 interactional features (see Appendix Table 2 for list) encoded and investigated to determine their diagnosis predicting power (using the same target variable as the non-interactional features). The following section explains the background, reasoning and hypothesis for choosing each variables and how it was encoded. The CHAT annotation was kept and used for the majority of interactional features and analysis of them. The quality of these annotations were sense checked by manually listening to a random selection of 25 transcripts (~5%) and sense checking that the annotations in the transcript were aligned with the recording.

A hierarchical grouping system of the interactional features was created which was based on a number of previous studies showing what interactional symptoms can be used

to diagnose AD and other types of dementia.[Svennevig and Lind 2016][Silvast, M. 1991][Itakura, H. 2001][Lesta, B et al. 2006][Cummings et al. 1985]. There were three hierarchical groupings that are referred to as ‘umbrella features’ where the sub-features within these groupings were different aspects of each umbrella feature. The three umbrella features are fillers, unintentional silence and conversation. The following is list of umbrella features, their sub-features and the reasoning and hypothesis behind choosing them to be investigated for this project followed by the variable name in brackets.

4.5.1 Umbrella Feature 1. Fillers

A filler is a sound or word that is spoken in conversation by one participant to signal to others a pause to think without giving the impression of having finished speaking. For example, “um”, “like”, “uh”, “you know!” and “actually.” As mentioned in the literature review, individuals with dementia of the Alzheimer’s type often experience, word retrieval problems (anomia) and typically as the disease progresses, the patient’s production of speech decreases (aphasia) and the use of empty phrases, speech automatisms increases, where fillers tend to be used to fill in the lexical gap [Svennevig and Lind 2016]. This is the reason that fillers were decided to be investigated as interactional features for this project.

A total of 10 different aspects of the filler term were investigated. It is important to note here that the CHAT annotation’s definition of a filler was used to locate the fillers throughout each session.

Total filler frequency was simply a count of the number of times a participant uttered a filler term. This was normalised by the total number of words spoken by the participant. (‘fillerCountPAR’)

POS tags post filler: Five features was encoded to investigate the different POS tags (Adjective, Adverb, Noun, Verb, Other) of the word that follows the filler term. The reason these features were created were to investigate the types of words AD patients were tending

to forget. It was hypothesised that participants with AD would be inclined to forget nouns more than any other POS. This hypothesis was constructed because, in [Jarrold et al. 2014], where acoustic and POS features were used to distinguish between 9 AD patients and 9 controls, confirmed that AD patients used more pronouns, verbs, and adjectives and fewer nouns.

A table (Appendix Table 5) of higher level POS tags was created so to not include granular POS tags. For example, the POS tag 'Verb', included all the different tenses of verb. ('postFillerPOS_Noun', 'postFillerPOS_Other', 'postFillerPOS_Verb', 'postFillerPOS_Adverb', 'postFillerPOS_Adjective')

A **'Filler Location' score** of each transcript was encoded. This takes into account the location of the filler in the utterance normalised by the total words of the transcript. It is hypothesised that AD patients will tend to forget a word sooner and have to pause and utter fillers earlier on in an utterances compared to non-AD patients. The location of the filler was defined by the index of the filler within the utterance. Therefore, if a filler occurred early in an utterance, it would have a low index value. I.e. If first word of utterance was a filler term, the index of this filler term is equal to zero. For each transcript, all index value were summed and this number was divided by the total number of fillers that were uttered by the participant, this value was a participants 'Filler Location' score. Participants with a low 'Filler Location' score were hypothesised to be diagnosed with AD as they would have uttered a high frequency of fillers (scores denominator) and a low value for the sum of filler term indexes (numerator). ('LocToFilRatio')

An **average cosine similarity score** (the same method was used to measure repetition of participants in the non-interactional features) of each word that occurred directly after the filler term was calculated for each transcript. If this variable was low this means that the all the words that the participant stated after the filler term were similar. This implies that thy participant was hesitating and forgetting words that were similar to each other. It was

hypothesised that this would represent a participant repeating the same mistakes. A high value for this variables would represent a participant that repeated the same mistakes, a known symptom of dementia. [Silvast, M. 1991] ('sim_array_post_fil')

A count of the following fillers were taken, based on the CHAT annotation, all normalised by the total number of words spoken by the participant:

A frequency count of the times a participant **laughed**. ('laughs')

A frequency count of the times a participant **sighed**. ('sighs')

Unintelligible words: A count of unintelligible words were encoded and used to signify the phenomena of mumbling, a common symptom of dementia [Lesta, B et al. 2006]. ('unintelligible')

Self-correction counts: A count of the number of times (based on CHAT notations) was also encoded. This variables could be used to represent self-repair, a representative when a person acknowledges that they made a mistake. If this variable negatively correlates with the AD diagnosis, it is hypothesised that this is because the mistake was made however, if it positively correlated with the target, this could be because a non-AD participant has acknowledged their mistake and corrected themselves. This variables fall under the filler umbrella feature because fillers and pauses tend to be used in the process of self-repair between the word that needs repairing and the speakers intended word. ('revisionCountsOnce')

4.5.2 Umbrella Feature 2. Unintentional Silence

This umbrella feature was used to represent when a moment of silence occurred or an utterance was interrupted or terminated unintentionally. This can happen for a number of reasons such as the participant loosing concentration or forgetting what they were supposed to be doing [Alzheimer's Society - What is Dementia? 2017]. This has similar reasoning for occurring as the above fillers such as aphasia and word retrieval [Silvast, M. 1991] however,

instead of trying to fill the lexical gap with semantically meaningless paraphrases or nouns or other speech automatism, the void in conversation is left as it is.

6 different aspect variables of this umbrella feature were encoded using the CHAT annotations and investigated. A count of the different **pause lengths** ('longPause', 'medPause', 'shortPause') was created. A variable to represent the **total length of pauses** uttered through out the session ('allPauseCount') was encoded by initializing a counter and adding a value of 1 to it when the participant uttered a short length pause, 2 when they uttered a medium length and 3 when a long length pause was uttered. All of the above pause based variables were normalised by the total number of pauses uttered by the participant. A count of when a participant **trailed off mid-utterance** and did not complete the utterance but was also not interrupted was also encoded. ('trailingOff')

Incomplete words: The number of incomplete words was also counted, however one the issues with this variables is that it may count a word as unfinished when it was actually the participants accents. E.g. If the participant said "travelin'" instead of "travelling", the programme would flag this as one incomplete word. This was normalised over the total number of words uttered by the participant. ('incompleteWordRatio')

4.5.3 Umbrella Feature 3. Conversation

This umbrella feature was created to represent phenomenon between the participant and the examiner that based on the conversation dominance. Dominance in every day conversation has been measured by the distribution between speakers of various interactional features, including topic control, interruptions and overlaps, and amount of speech. [Itakura, H. 2001] The reason this area of interaction was investigated for this project was due to [Silvast, M. 1991], an explorational study carried out between speech therapists and patients of aphasia. It's results showed that during the conversation therapists had a role which was manifested in their frequent use of requests for information and

clarification. Aphasics had more speech time but were in a responsive role in the conversation. Aphasia is a common trait of AD patients, which is shown in [Cummings et al. 1985] where a speech and language assessment in 30 patients with dementia of the Alzheimer type and in 70 normal controls revealed that all AD patients were aphasic. This is linked to the decline in AD patients ability to concentrate, making it difficult for the patient to complete tasks and to follow conversation. [Alzheimer's Society - What is Dementia? 2017] Therefore, it is hypothesised in this project, if the participant has AD, the examiner, similar to [Silvast, M. 1991], will need to encourage the participant for more information about the picture throughout their descriptive task. This will result in shorter participant answers to examiner questions, a higher turn-taking rate between the two parties and the examiner having a higher presence throughout the session for participants with AD and other dementias than that of the control.

A total of 7 variables were encoded under this umbrella feature. The following is a list of variables and the algorithm used to encode these interactional phenomenon between the participant and the examiner:

Backchannels: A ratio of the number of times the examiner used a backchannel, such as “Mhm” to signify approval or a suggestion that they are expecting more information from the participant, was calculated by counting the number of backchannels (based on the CHAT annotations) and divided by the number of words spoken by the participant. Therefore, it is suggested that a participant with a high value for this variable would have AD. ('backChannelmhm')

Questions: A count of the both the examiner and the participants questions were taken and normalised over their total number of utterances respectively. ('quesCountPAR2', 'quesCountINV2') as well as a count of the times the examiner was interrupted by a participant by an utterance in the shape of a question ('Interruption_Q'). The final question variable was created based on the CHAT annotation.

Answer Length: The average answer length per transcript was calculated first counting the number of words the participant stated directly after the examiner's utterances ended in a question mark and up until the next utterance commenced. This value was then divided by a count of the examiners question marks in order to find the average answer length per transcript. ('avgAnsWordLen2')

Turn Count Ratio: A total count of turns (turnCountINV/turnCountPAR) were also counted within each session. Therefore, the higher the value for this variable, the more utterances the examiner had through out the session. ('turnCountRatio')

The number of times a participant uttered "I don't know" or synonyms of that phrase was counted and normalised. This can be used to represent a lack of clarity of the task at hand or misunderstanding between the participant and the examiner. ('dunnoCount')

4.6 Non-Interactional Features Selection

A total of 34 non-interactional features were originally encoded (based on [Fraser et al. 2016]'s non-acoustic top predicting variables) for this project. The top predicting variables, regarding the classification of the diagnosis, were selected from these 34 features. The aimed outcome was a classifier trained on the TalkBank data using these encoded non-interactional features as variables with a resulting similar accuracy score to that of the classifier using similar features and the same data in [Fraser et al. 2016] to create a baseline to evaluate the interactional features against. In order to do this, it was decided to keep as many variables built as possible, as they are already proven to optimal predictors, and base the feature selection method on the removal of variables that were correlated with each other. This was due to the fact that the phenomena that was being represented was already encoded in some other variables and leaving the overfitting features in can lead to the model overfitting resulting in higher accuracy scores that are not representative of the classifier.

The first step was to investigate the correlation between each of the variables. This was done by creating a correlation matrix (Appendix Fig. 1) with all 34 linguistic features. As anticipated and seen from the correlation matrix below, variables that were built on counts of the letters or words correlated. For example, the average word length of the session (1: 'Avg_Word_Len') correlates with the number of letters within the script (13: 'Script_Letter_Ct') so it was decided to remove the latter. The total count of verbs (23: 'Verb_Ct_Total') naturally correlated with the break down of all the different verb tenses (17: 'Verb_Ct_3sing', 18: 'Verb_Ct_Base', 19: 'Verb_Ct_Gerund', 20: 'Verb_Ct_Past', 21: 'Verb_Ct_Past_P', 22: 'Verb_Ct_Past_T', 24: 'Verb_Ct_non3sing') as well as the participant's number of ProNouns (11: 'ProNoun_Ct') and the total number of utterances spoken by the participant (9: 'PAR_UTT_ct'). Since it correlated with so many variables it was thought best to, not only run a correlation matrix with it removed and but to also investigate the correlation matrix with the total count of verbs remaining instead of the break down of different verbs as well as the different combination of the verb break downs. Any variables where there only difference between them and another variable was the normalisation step, were removed. All the parser based variables remained as they didn't correlated with themselves or any other variables. 11 variables in total were removed in this step so 23 remained. The correlation matrix was built again to confirm there was less correlation between the variables (Appendix Fig. 2). The remaining 23 variables after the feature selection (features in bold font in above legend) were decided to be the top non-interactional features to train the classifier on.

4.7 Classifier Build

Four different classifiers were trained, each on four different cuts of the data labelling's of the target variable. The four classifiers were decision tree, logistic regression, Gaussian Naive Bayes and k-nearest neighbour, which were all built by utilizing the Python library

'sklearn'. The four different cuts of the data and labelling's of the data (D1, D2, D3, D4) can be seen in Appendix Table 4.

The decision tree classifier split the data on the best split with its random state parameter being set to 0. The logistic regression classifier used an L2 regularizer to avoid overfitting. Both the Gaussian Naive Bayes and k-nearest neighbour classifier used the default 'sklearn' parameters.

Due to the relatively small data set of (550 entries, each representing one participant), 10-fold cross validation was used to train and test the data on each classifier in order to avoid overfitting and ensure that the test set was representative of the entire data set. Regarding the target labels, the data was unbalance; there was a large number of 'Probable AD' and 'Control' participants and a low count of the other targets. To ensure a fair performance metric, precision, recall and F1 scores were all calculated for each fold for each classifier. The average of the 10 folds was then taken as the final performance metric (see Appendix Table 3). The different labelling's of the target variables and cuts of the of the data that were mentioned above were also compared.

[Fraser et al. 2016] built their model using logistic regression and the same cut of data, and target labels as 'Dataset 4' (D4) in Table 2 below. This merged both severities ('Possible' & 'Probable') of AD as 'All AD' and compared them against the control only. All other dementia causes were ignored. As can be seen from Table 1 below, when the LR classifier was ran on the D4 dataset for this project, the accuracy scores were almost 78%, with both a recall and F1 score of approximately 73% (highlighted in bold in Table 1 below). These accuracy results are similar to that of [Fraser et al. 2016] (~80%) where they applied the same classifier to the same cut of data. It was expected that this projects performance metrics would be slightly lower than those of [Fraser et al. 2016] since the acoustic features were not encoded and so less features were used.

It has been shown in [Fraser et al. 2016] what the optimal non-interactional features are to predict the different diagnosis of AD. This project re-encoded these features (NIF's) and removed variables that correlated amongst each other. Due to the similarity in performance metrics of the LR classifier for D4 in this project and those of [Fraser et al. 2016] where the same data, target labels and classifier method were used, these remaining variables (Top NIF's) can be used as a baseline to evaluate the performance of the interactional features.

X = Top NIF's (Count 23)		NIF Only Performance (10-fold CV)		
y	Classifier	Precision	Recall	F1 Score
D1:	DT	0.5233	0.4727	0.4917
	LR	0.5881	0.6200	0.5948
	GNB	0.5999	0.1236	0.1682
	KNN	0.5277	0.5145	0.5046
D2:	DT	0.6197	0.5670	0.5810
	LR	0.6991	0.6831	0.6821
	GNB	0.7587	0.2306	0.3096
	KNN	0.6136	0.5789	0.5838
D3:	DT	0.5160	0.4636	0.4816
	LR	0.5889	0.6200	0.5953
	GNB	0.5983	0.1309	0.1745
	KNN	0.5346	0.5109	0.5044
D4:	DT	0.7117	0.6293	0.6517
	LR	0.7750	0.7233	0.7381
	GNB	0.7488	0.6914	0.7094
	KNN	0.7102	0.6231	0.6455

Table 1. Count of data for the different labelling's of the target variable.

	Control	Possible AD	Probable AD	MCI	Vascular	Memory
D1:	241	21	237	43	5	3
D2:	241	21	237	0	0	0
D3:	241	21	237	51 ('Other')		
D4:s	241	258 ('All AD')		0	0	0

Table 2. Count of data for the different labelling's of the target variable

The logistic regression classifier performed the best across all the different cuts and labelling's of the data, in particular for D2 & D4. However, this is because there is a low frequency of other 'Non-AD' or 'Control' labels, hence the relative low recall and F1 scores. Instead of removing the less common target entries, it was decided to base the analysis for this project on D4, where all 'Non-AD' or 'Control' targets were removed and both severities ('Possible' and 'Probable') of AD were merged to return a target label of 'All AD'.

4.8 Interactional Features Selection

A total of 26 interactional features were encoded originally. The feature selection process for these features was done by ranking all 26 encoded features based on their ANOVA F-value between label/feature for the classification task of predicting the diagnosis. Using Python's SelectKBest (sklearn) function where the number of features that are being ranked, k, is inputted as a parameter. In this case, k is equal to 26. It was decided that the top interactional features (Top IF's) would be selected based on their ANOVA F-value score (predicting weight) which is plotted in Appendix Figure 3. All features with a non-zero ANOVA F-value score were selected as the Top IF's therefore the only feature that was removed was 'Interruption_Q' which was encoded to represent if the participant interrupted either themselves or the examiner in the form of a question. Therefore, Top IF's were made up of 23 features (see Appendix Table 2).

From Appendix Figure 3, it can be seen that the top two predicting features based on ANOVA F-Value all appeared to fall under the 'Conversational' umbrella feature that was created. The ratio of utterance starts between the participant and examiner (Turn-Taking Ratio) and a normalised count of questions asked by the invigilator, which can both be said to represent the examiner needing clarify with the participant that they comprehend the task at hand or suggests that the examiner was required to encourage them for more information. The fact that the third highest predicting variable (Incomplete Word Count of the participant) suggests that this could be the cause for the need for the examiner to have to ask more questions, and so have a higher utterance start rate. The combination of these three features could represent the amount of involvement that the examiner needed to manifest throughout the task. Further exploration of the interactional features was carried out by investigating the correlation between the target variable and the each of the Top IF's.

4.9 Diagnosis Correlation Analysis of Interactional Feature

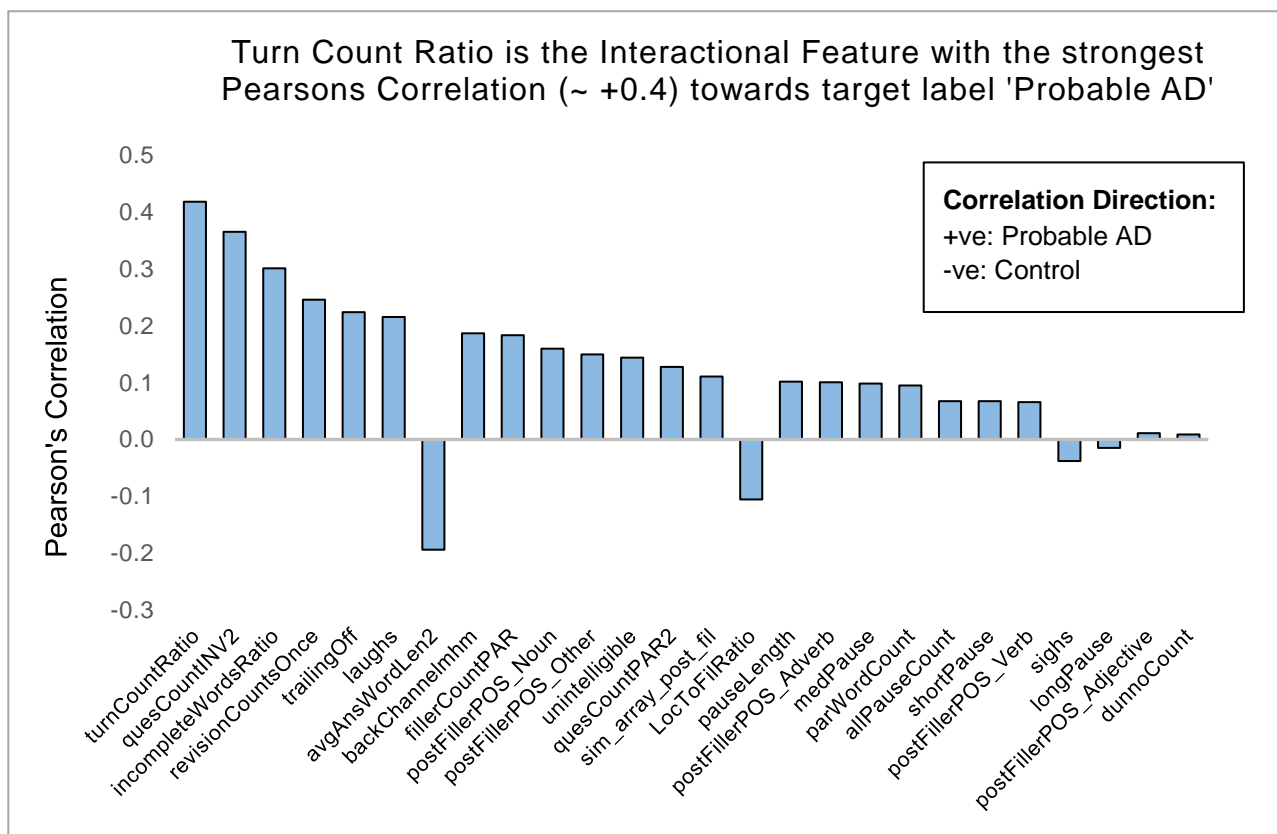


Fig 1. Correlation between target variables ('Probable AD') and the top interactional features.

The aim of this section⁶ was to investigate how each of the interactional variables correlate with the target variables (diagnosis). The also acts as a sense check that the variables were in encoded as expected. The strength and direction (positive or negative) of each variable's correlation with the target was investigated. To ensure a representative output, it was decided to compare the control group against the group with the highest level of dementia (Probable AD) using Pearson's correlation where if a variable has a positive correlation, it's value increases in the direction of 'Probable AD' and if it has a negative correlation, it's value increases in the direction of the control target.

⁶All data in the section refers to Fig. 1 (above) unless stated otherwise.

The two variables that have the strongest positive correlation (TurnCountRatio and Examiner's Question Frequency ('quesCountINV2')) both suggest that the examiner needed to manifest a stronger presence throughout the session, encouraging the participant for more information or having to clarify that the participant understands. This theory is strengthened by the fact that 'backChannelhmh' which represents is the normalised count of backchannels used by the examiner also positively correlates with the diagnosis, and overlaps with the reasoning behind the variable with strongest negative correlation (average answer word length). It suggests that the longer the length of the participants sentence, the more likely they are to not have 'Probable AD'.

The number of fillers (normalised by being divided by the total words) uttered by the participant positively correlates with the participant having AD. This is also said for the 'revision count' variable, which can be said to roughly represent self-repair. Since it positively correlates with the diagnosis, it suggests that the AD patients recognise when they say the unintended word that needs to be repaired and attempt to repair it.

The slightly negative correlation of the 'locToFillerRatio' variable and the target suggests that there is some relationship between the two although not as strong as the previous variables mentioned. It suggests that the earlier on in a sentence (lower locToFillerRatio value) that the filler occurs, the more likely the participant is to have AD.

The 'Probable AD' diagnosis correlate's with participants uttering a filler, which can be used as a sign of hesitation, before nouns and adverbs, where as there is little correlation with the target variables and a patient utters a filler before an adjective, which was expected as it has be shown in [Jarrold et al. 2014] that AD patients utter few nouns that other POS tags however, it has now been shown that this is also through for the POS tags that occur after a filler.

5. Evaluation Methods

To evaluate the performance of interactional features being used to computationally diagnose Alzheimer's types and other dementia, the interactional features would be merged with variables that were already known to be strong predictors of the diagnosis (ie. the top Non-IF's), to see if any of the interactional features had a higher predicting weight (ANOVA F-Value) than any of the Non-IF's. If any of the interactional features outperformed the non-interactional features regarding their predicting weight score, it can be said that the combination of both interactional and non-interactional features predict better than that of just non-interactional features.

The null hypothesis was that the top predicting variables regarding the computational diagnosis of Alzheimer types are based on non-interactional, linguistic based phenomenon. These are represented by the top 23 non-interactional features (Top Non-IF's) that were created based on the non-interactional features in [Fraser et al. 2016]. The alternative hypothesis is that the addition of other variables, encoded based on the interactions of the patient, are better predictors and can be used to improve the classification of diagnosis.

This is the outcome if any of the Top IF's fall into the top 23 variables based on predictive weight when both the Top Non-IF's and Top IF's predictive weights are compared. The 'new' top 23 variables based on predictive power will then be used to train the same classifier and same data that was built earlier using the top Non-IF's. If there is an improvement in performance scores, it cements the statement that the combination of both interactional and non-interactional features predict better than that of just non-interactional features regarding the computational diagnosis of Alzheimer's disease.

6. Results

The final top 23 predicting features were made up of 52.2% interactional and 47.8% non-interactional. (See Appendix Table 6) This combination of features improved the originally all non-interactional classifier accuracies by 4.96%, with an accuracy score of 82.46%, recall at 78.16% and an F1 Score of 79.36% (See below Table 3 or Appendix Table 3).

X = Top NIF's (Count 23)	Classifier	'NIF Only' Performance (10-fold CV)			'NIF & IF' Performance (10-fold CV)			Absolute difference (from 'NIF only' to 'NIF & IF')		
		Precision	Recall	F1 Score	Precision	Recall	F1 Score	Precision	Recall	F1 Score
y	DT	0.7117	0.6293	0.6517	0.7414	0.6836	0.7006	2.97%	5.43%	4.90%
	LR	0.7750	0.7233	0.7381	0.8246	0.7816	0.7936	4.96%	5.82%	5.54%
	GNB	0.7488	0.6914	0.7094	0.8023	0.7333	0.7496	5.35%	4.19%	4.03%
	KNN	0.7102	0.6231	0.6455	0.6916	0.5791	0.6043	-1.86%	-4.39%	-4.12%

Table 3. Performance scores of all classifiers dataset 4: 'All AD' only against Ctrl

In terms of the interactional variables 'Umbrella Features' and the non-interactional variables 'Factor Groupings', the final top 23 variables consisted of 17.4%: 'Fillers', 30.4%: 'Semantic Impairment', 17.4%: 'Information Impairment', 13.0%: 'Conversation' and 17.4%: 'Unintentional Silence'. (See Appendix Fig 6)

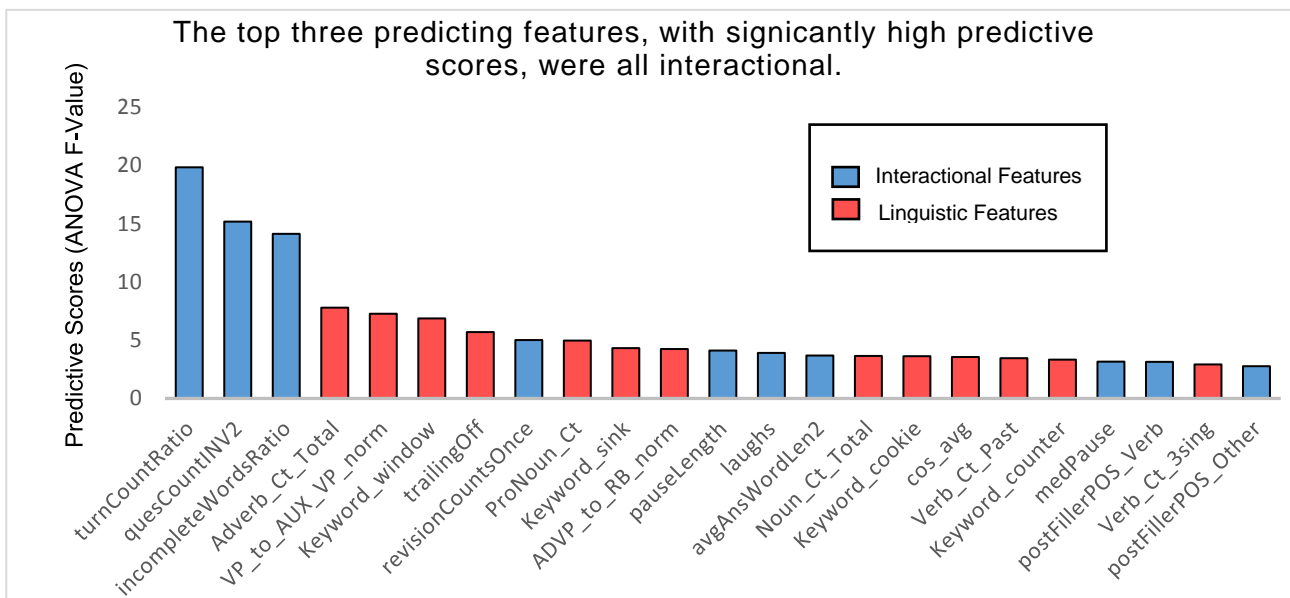


Fig 2. Predictive weights of 'New Top 23' Features (data from Appendix Table 6)

There are three variables that have a significantly higher predicting scores (ANOVA F-Value) than the remaining features, which are, Turn-Count Ratio (TCR), 'quesCountINV2'

(Examiner Question Count - EQC) and Incomplete Word Ratio (IWR). The combination of these three variables having a significantly higher predictive scores suggest that the phenomenon of the examiner being required to manifest a role of higher dialogue and presence can be represented by encoding these features and used to aid the computational diagnosis of AD.

6.1 Discussion

It is clear from the above results that the combination of features that represent both linguistic and interactional phenomenon perform better at computationally diagnosing AD than using features that are solely based on linguistics. This is due to the fact that, as hypothesized earlier on the in project, by adding in an extra layer of known symptoms to the classifier, it will improve predictions and classification of AD.

Regarding the target (diagnosis) predictive weight score the two highest ranking variables, 'Turn Count Ratio' (defined as 'examiner utterance starts'/'participant utterance starts') and a count of the questions the examiner stated (EQC), were both calculated based on the amount of involvement and presence the examiner was required to manifest throughout the session, with predictive scores of 19.85 and 15.18 respectively. The third highest predicting variable, with a predictive score of 14.14, represents a count of incomplete words (IWR) uttered by the participant throughout the session. When compared with the correlation analysis (interactional features against the diagnosis variable) carried out previously, it can be seen that these three top predicting features all positively correlate with the diagnosis of 'Probable AD'. A suggested reasoning to this is based on the studies of [Silvast, M. 1991], which shows that the examiner will manifest a role that will assist an aphasia patient when needed. This is common in the field of speech and language therapy interactions, where the examiner takes a scaffolding role in the interaction. In this case it can be seen that the examiner has a need to take on a role with a higher level of presence

(higher TCR) throughout the session due to the fact that the participant is unclear (higher IWR) leading to the examiner requiring to ask questions (higher EQC & TCR) in order to clarify what the participant is intending to express and to encourage the participant to keep going. The combination of these three variables having a significantly higher predictive score may be used to suggest that the phenomenon of the examiner being required to manifest a role of higher dialogue and presence can be represented by encoding these features and similar features to these, in order to aid the computational diagnosis of AD. This leads to the hypothesis that, if encoded correctly, conversation and dialogue analysis may play a significant role in computationally classifying Alzheimer's type. This could be proven by investigating dialogue-based features.

6.2 Conclusion

Since interactional features replaced over half of the originally all non-interactional top predicting features regarding diagnosis predicting weight, as well as causing an improvement in the classifier's accuracy, one can conclude that encoding interactional features, in particular dialogue based features that represent the amount of involvement the invigilator is required to manifest throughout the session, in addition to non-interactional features, can assist in computationally classifying Alzheimer's disease.

6.3 Further Work

The results suggest that dialogue plays a key role in the classification of AD therefore further investigation would include exploration of variables created based on the dialogue between the examiner and the participant. This would include running the data through a dialogue act tagger and encoding variables that represent deeper examiner-participant interactions and carrying out an analysis similar to the above, including a correlation analysis with the target variable.

It would have also been beneficial to investigate a different data set, other than recordings of the “Boston Cookie Theft” task as other psychological tasks may include a stronger level of interactions between the examiner and the patients. It was found that, although there are a number of other tests available (Sentence⁷, Fluency⁸, Recall⁹) on DementiaBank, there is no associated control group. Therefore, in order to ensure a fair analysis of these data, a proxy control group would have to be created.

It would be interesting to carry out a sentiment analysis on a dataset that expresses patients emotion and encode a sentiment based variable for AD classification, a linguistic approach to work of [Lopez-de-Ipiña et al. 2015], where emotional temperature is calculated using speech signal as an AD classification feature. Further work could involve investigating how patient-examiner interactions effect the sentiment of the patient or if the sentiment change is the variable that causes a change in the patient-examiner interactions. This could be used to improve the efficiency of certain patient-examiner framework used in general speech language sessions.

Self-repair plays a significant role when in the process of word-searching (anomia). This is a key symptom of dementia due to aphasia. The self-correction variable used in this project appeared in the top ten AD predicting features. It was also felt however, that the definition of the self-correction (based on CHAT annotation of ‘self-restarts’) variable could be encoded in a more robust manner by running the data through a high end self-repair detector such as ‘STIR’, an automatic incremental self-repair detection system developed on the Switchboard corpus of telephone [Hough 2014]. It is felt that a number of different

⁷The examiner gives examples of words and instructs the participant to to use the word in a sentence.

⁸Examiner instructs the participant to name as many animals as they can within one minute.

⁹The examiner tells a store and instructs the participant to retell everything that they can remember from the story.

variables could be built based on self-repair including finding the cosine difference between word that needed repairing and intended words. This could be said to measure the mistake made.

This project focuses on diagnosis Alzheimer's disease, and even though this is the most common form of dementia, it would be interesting to investigate what features predict other NON-AD dementia (MCI, Vascular, Memory). As was seen earlier on in the project, this dataset is very unbalanced regarding non-AD target labels. Therefore, one could use an anomaly detection approach as opposed to classification and investigate the features that aid the prediction.

7. Bibliography

1. Fraser et al. 2016: Fraser, K.C., Meltzer, J.A. and Rudzicz, F., Linguistic features identify Alzheimer's disease in narrative speech., 2016
2. Alzheimer's Society - What is Dementia? 2017: Alzheimer's Society, Alzheimer's Society - What's is Dementia?, 2017,
https://www.alzheimers.org.uk/info/20007/types_of_dementia/1/what_is_dementia
3. Alzheimer's Society - Facts for the media 2017: Alzheimer's Society, Facts for the media - Alzheimer's Society., 2017
4. Svennevig and Lind 2016: Svennevig, J. and Lind, M., Dementia, interaction, and bilingualism: An exploratory case study., 2016
5. Croisile et al. 1996: Croisile, B., Ska, B., Brabant, M.J., Duchene, A., Lepage, Y., Aimard, G. and Trillet, M., Comparative study of oral and written picture description in patients with Alzheimer's disease., 1996
6. Guinn and Habash 2012: Guinn, C.I. and Habash, A., Language Analysis of Speakers with Dementia of the Alzheimer's Type., 2012
7. Bucks et al. 2000: Bucks, R.S., Singh, S., Cuerden, J.M. and Wilcock, G.K., Analysis of spontaneous, conversational speech in dementia of Alzheimer type: Evaluation of an objective technique for analysing lexical performance., 2000
8. Alzheimer's Society - Behaviour Changes 2017: Alzheimer's Society, Alzheimer's Society - Behaviour Changes, 2017,
https://www.alzheimers.org.uk/info/20064/symptoms/87/behaviour_changes
9. Horn et al. 2009: Horn, J.F., Habert, M.O., Kas, A., Malek, Z., Maksud, P., Lacomblez, L., Giron, A. and Fertil, B., Differential automatic diagnosis between Alzheimer's disease and frontotemporal dementia based on perfusion SPECT images., 2009
10. Klöppel et al. 2008: Klöppel, S., Stonnington, C.M., Chu, C., Draganski, B., Scahill, R.I., Rohrer, J.D., Fox, N.C., Jack Jr, C.R., Ashburner, J. and Frackowiak, R.S., Automatic classification of MR scans in Alzheimer's disease., 2008
11. Lopez-de-Ipiña et al. 2015: Lopez-de-Ipiña, K., Alonso, J.B., Solé-Casals, J., Barroso, N., Henriquez, P., Faundez-Zanuy, M., Travieso, C.M., Ecay-Torres, M.,

- Martinez-Lage, P. and Eguiraun, H., On automatic diagnosis of Alzheimer's disease based on spontaneous speech analysis and emotional temperature., 2015
12. DementiaBank 1994: Becker, J. T., Boller, F., Lopez, O. L., Saxton, J., & McGonigle, K. L., The natural history of Alzheimer's disease: description of study cohort and accuracy of diagnosis., 1994, <http://dementia.talkbank.org/>
 13. Goodglass and Kaplan 1983: Goodglass H, Kaplan E, The Boston Diagnostic Aphasia Examination., 1983
 14. Silvest, M. 1991: Silvest, M., Aphasia therapy dialogues., 1991
 15. Fergadiotis et al. 2011: Fergadiotis, G. and Wright, H.H., , Lexical diversity for adults with and without aphasia across discourse elicitation tasks., 2011
 16. Szatloczki et al. 2015: Szatloczki, G., Hoffmann, I., Vincze, V., Kalman, J. and Pakaski, M., Speaking in Alzheimer's disease, is that an early sign? Importance of changes in language abilities in Alzheimer's disease., 2015
 17. TalkBank 2007: MacWhinney, B., The TalkBank Project., 2007
 18. Becker et al. 1994: Becker JT, Boiler F, Lopez OL, Saxton J, McGonigle KL, The natural history of Alzheimer's disease: Description of study cohort and accuracy of diagnosis., 1994
 19. Carozza 2015: Carozza, L.S. ed., Communication and Aging: Creative Approaches to Improving the Quality of Life., 2015
 20. Giles et al. 1996: Giles, E., Patterson, K. and Hodges, J.R., Performance on the Boston Cookie Theft picture description task in patients with early dementia of the Alzheimer's type: missing information., 1996
 21. Bird 2000: Bird, H., Ralph, M.A.L., Patterson, K. and Hodges, J.R., The rise and fall of frequency and image ability: Noun and verb production in semantic dementia., 2000
 22. MFMER 2016: Mayo Foundation for Medical Education and Research (MFMER), Mild Cognitive Impairment (MCI), 2016, <http://www.mayoclinic.org/diseases-conditions/mild-cognitive-impairment/home/ovc-20206082>
 23. Alzheimer's Society - Vascular Dementia 2017: Alzheimer's Society, Vascular Dementia, 2017, <http://www.alzheimersresearchuk.org/about-dementia/types-of-dementia/vascular-dementia/about/>
 24. MacWhinney, B. 2000: MacWhinney, B., The CHILDES Project: Tools for analyzing talk, 3rd edition., 2000
 25. Honoré 1979: Honoré, A., Some simple measures of richness of vocabulary., 1979

26. Bloom et al. 1980: Bloom, L., Lifter, K. and Hafitz, J., Semantics of verbs and the development of verb inflection in child language., 1980
27. Itakura, H. 2001: Itakura, H., Describing conversational dominance., 2001
28. Lesta, B et al. 2006: Lesta, B. and Petocz, P., Familiar group singing: Addressing mood and social behaviour of residents with dementia displaying sundowning., 2006
29. Cummings et al. 1985: Cummings, J.L., Benson, D.F., Hill, M.A. and Read, S., Aphasia in dementia of the Alzheimer type., 1985
30. Jarrold et al. 2014: Jarrold, W., Peintner, B., Wilkins, D., Vergryi, D., Richey, C., Gorno-Tempini, M.L. and Ogar, J., Aided diagnosis of dementia type through computer-based analysis of spontaneous speech., 2014
31. Hough 2014: Hough, J., (Doctoral dissertation, Queen Mary University of London).
32. Dingemanse et al. 2013: Dingemanse, M., Torreira, F. and Enfield, N.J., 2013. Is “Huh?” a universal word? Conversational infrastructure and the convergent evolution of linguistic items. PloS one, 8(11), p.e78273. 2013

8. Appendix

Fig 1: Correlation Matrix of All NI Features (Pre-FS)

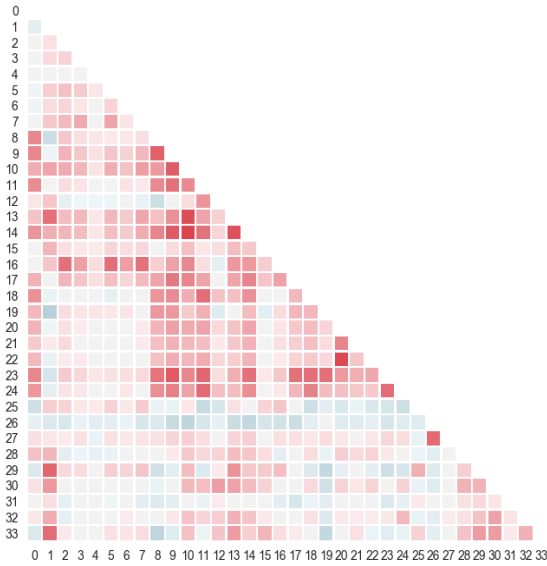


Fig 2: Correlation Matrix of Remaining NI Feature's (Post-FS)

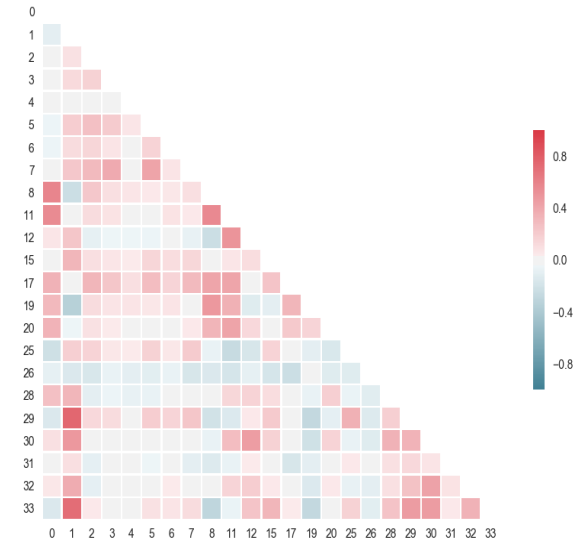


Table 1. Fig. 1 and Fig. 2 legend & list of NI variables encoded

0: 'Adverb_Ct_Total'	12: 'ProNoun_Noun_Ratio'	24: 'Verb_Ct_non3sing'
1: 'Avg_Word_Len'	13: 'Script_Letter_Ct'	25: 'cos_avg'
2: 'Keyword_cookie'	14: 'Script_Word_Ct'	26: 'wordsOnceOverTotal'
3: 'Keyword_counter'	15: 'Subj_IU_Total'	27: 'wordsUsedOnce'
4: 'Keyword_curtain'	16: 'Total_Keywords_Mentioned'	28: 'ADVP_to_RB_norm'
5: 'Keyword_sink'	17: 'Verb_Ct_3sing'	29: 'NP_to_DT_NN_norm'
6: 'Keyword_stool'	18: 'Verb_Ct_Base'	30: 'NP_to_PRP_norm'
7: 'Keyword_window'	19: 'Verb_Ct_Gerund'	31: 'ROOT_to_FRAG_norm'
8: 'Noun_Ct_Total'	20: 'Verb_Ct_Past'	32: 'VP_VBG_norm'
9: 'PAR_UTT_ct'	21: 'Verb_Ct_Past_P'	33: 'VP_to_AUX_VP_norm'
10: 'PAR_Word_Ct'	22: 'Verb_Ct_Past_T'	
11: 'ProNoun_Ct'	23: 'Verb_Ct_Total'	

*PAR: Participant, Ct: Count, UTT: Utterance, IU: Information Unit

Table 2. List of all interactional features encoded.

0: Interruption_Q	9: medPause	18: trailingOff
1: LocToFilRatio	10: pauseLength	19: turnCountRatio
2: allPauseCount	11: postFillerPOS_Noun	20: unintelligible
3: backChannelmhm	12: postFillerPOS_Other	21: avgAnsWordLen2
4: dunnoCount	13: postFillerPOS_Verb	22: quesCountINV2
5: fillerCountPAR	14: revisionCountsOnce	23: quesCountPAR2
6: incompleteWordsRatio	15: shortPause	24: postFillerPOS_Adjective
7: laughs	16: sighs	25: postFillerPOS_Adverb
8: longPause	17: sim_array_post_fil	

*'0: Interruption_Q' was the only feature removed during the feature selection process.

Fig. 3 Predictive Scores (X-axis) vs. All IF's (Y- axis)

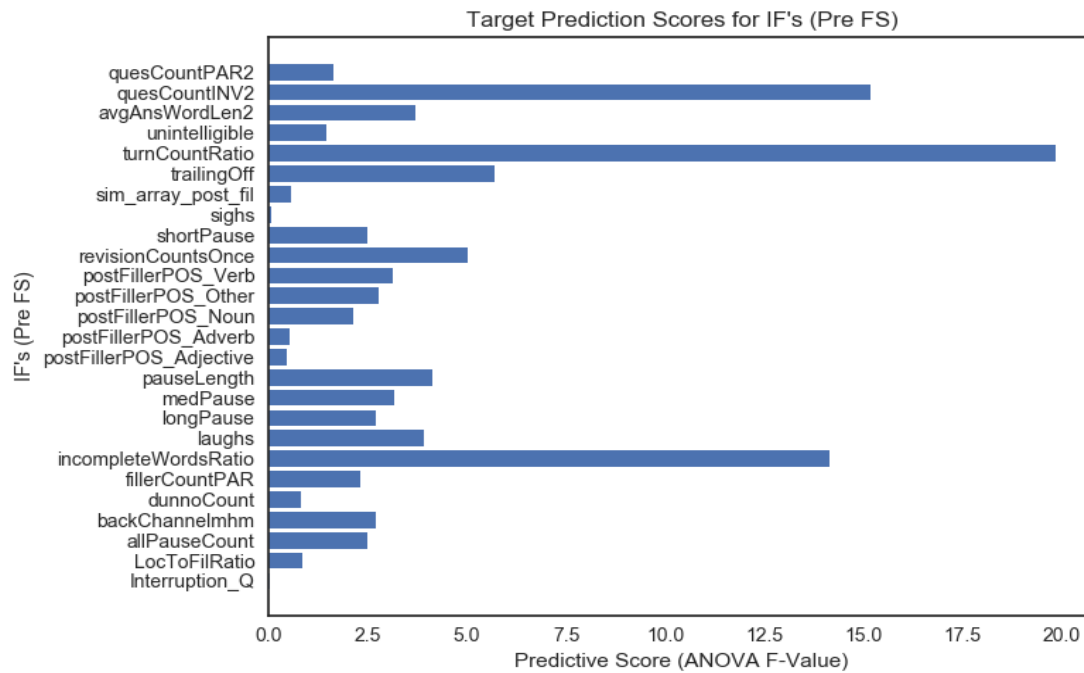


Table 3. Performance scores of all classifiers on different cuts of the data.

X = Top NIF's (Count 23)		'NIF Only' Performance (10-fold CV)			'NIF & IF' Performance (10-fold CV)			Absolute difference (from 'NIF only' to 'NIF & IF')		
y	Classifier	Precision	Recall	F1 Score	Precision	Recall	F1 Score	Precision	Recall	F1 Score
D1:	DT	0.5233	0.4727	0.4917	0.6275	0.5364	0.5645	10.42%	6.36%	7.27%
	LR	0.5881	0.6200	0.5948	0.6367	0.6818	0.6506	4.87%	6.18%	5.58%
	GNB	0.5999	0.1236	0.1682	0.6568	0.3436	0.4298	5.69%	22.00%	26.17%
	KNN	0.5277	0.5145	0.5046	0.5620	0.5145	0.5110	3.43%	0.00%	0.64%
D2:	DT	0.6197	0.5670	0.5810	0.6680	0.5909	0.6157	4.83%	2.39%	3.47%
	LR	0.6991	0.6831	0.6821	0.7626	0.7573	0.7523	6.35%	7.42%	7.02%
	GNB	0.7587	0.2306	0.3096	0.7454	0.4627	0.5377	-1.33%	23.21%	22.81%
	KNN	0.6136	0.5789	0.5838	0.6475	0.5770	0.5904	3.40%	-0.19%	0.66%
D3:	DT	0.5160	0.4636	0.4816	0.6014	0.5200	0.5432	8.54%	5.64%	6.15%
	LR	0.5889	0.6200	0.5953	0.6344	0.6818	0.6496	4.55%	6.18%	5.42%
	GNB	0.5983	0.1309	0.1745	0.6538	0.3873	0.4605	5.55%	25.64%	28.60%
	KNN	0.5346	0.5109	0.5044	0.5571	0.5145	0.5124	2.25%	0.36%	0.80%
D4:	DT	0.7117	0.6293	0.6517	0.7414	0.6836	0.7006	2.97%	5.43%	4.90%
	LR	0.7750	0.7233	0.7381	0.8246	0.7816	0.7936	4.96%	5.82%	5.54%
	GNB	0.7488	0.6914	0.7094	0.8023	0.7333	0.7496	5.35%	4.19%	4.03%
	KNN	0.7102	0.6231	0.6455	0.6916	0.5791	0.6043	-1.86%	-4.39%	-4.12%

Table 4. Count of data for the different labellings of the target variable for Appendix Table 3 (above).

	Control	Possible AD	Probable AD	MCI	Vascular	Memory
D1:	241	21	237	43	5	3
D2:	241	21	237	0	0	0
D3:	241	21	237	51 ('Other')		
D4:	241	258 ('All AD')		0	0	0

Fig. 4 Predictive Scores (X-axis) vs. Top IF's and NIF's (Y- axis)

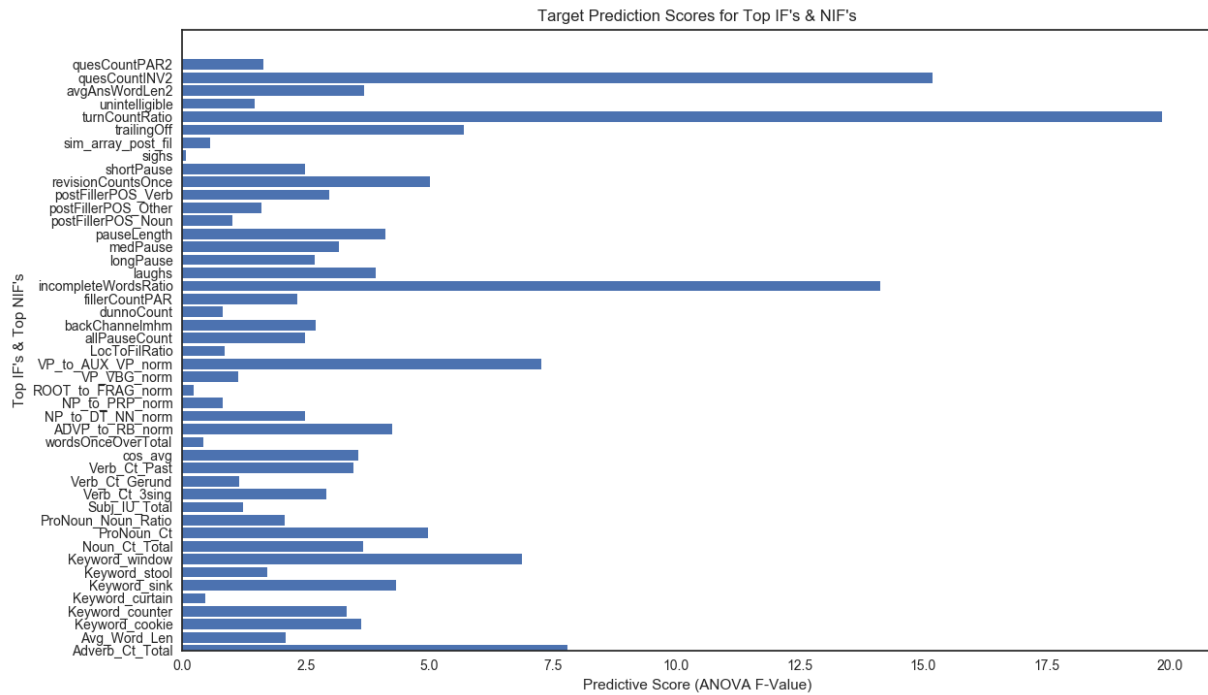


Table 5. POS Tags - High Level Groupings

POS Tag	Description	Grouping
CC	Coordinating conjunction	Other
CD	Cardinal number	Other
DT	Determiner	Other
EX	Existential there	Other
FW	Foreign word	Other
IN	Preposition or subordinating conjunction	Other
JJ	Adjective	Adjective
JJR	Adjective, comparative	Adjective
JJS	Adjective, superlative	Adjective
LS	List item marker	Other
MD	Modal	Other
NN	Noun, singular or mass	Noun
NNS	Noun, plural	Noun
NNP	Proper noun, singular	Noun
NNPS	Proper noun, plural	Noun
PDT	Predeterminer	Other
POS	Possessive ending	Other
PRP	Personal pronoun	Other
PRP\$	Possessive pronoun	Other
RB	Adverb	Adverb
RBR	Adverb, comparative	Adverb
RBS	Adverb, superlative	Adverb
RP	Particle	Other
SYM	Symbol	Other
TO	to	Other
UH	Interjection	Other
VB	Verb, base form	Verb
VBD	Verb, past tense	Verb
VBG	Verb, gerund or present participle	Verb
VBN	Verb, past participle	Verb
VBP	Verb, non-3rd person singular present	Verb
VBZ	Verb, 3rd person singular present	Verb
WDT	Wh-determiner	Other
WP	Wh-pronoun	Other
WP\$	Possessive wh-pronoun	Other
WRB	Wh-adverb	Other

Table 6. Final Top 23 Features and their Predictive Weights

Rank	Predictive Scores (ANOVA F-Value)	Variable	Feature Type	Factor/Umbrella
1	19.84673728	turnCountRatio	IF	Conversation
2	15.1849209	quesCountINV2	IF	Conversation
3	14.13549536	incompleteWordsRatio	IF	Unintentional Silence
4	7.796872585	Adverb_Ct_Total	NIF	Semantic
5	7.275023488	VP_to_AUX_VP_norm	NIF	Syntactic
6	6.873187562	Keyword_window	NIF	Information
7	5.700162083	trailingOff	NIF	Unintentional Silence
8	5.015922563	Self-Repair (revision count)	IF	Filler
9	4.969879595	ProNoun_Ct	NIF	Semantic
10	4.323192073	Keyword_sink	NIF	Information
11	4.245024757	ADVP_to_RB_norm	NIF	Semantic
12	4.114834213	pauseLength	IF	Unintentional Silence
13	3.912345997	laughs	IF	Filler
14	3.687530286	avgAnsWordLen2	IF	Conversation
15	3.653043837	Noun_Ct_Total	NIF	Semantic
16	3.629273321	Keyword_cookie	NIF	Information
17	3.562162496	cos_avg	NIF	Semantic
18	3.458298487	Verb_Ct_Past	NIF	Semantic
19	3.33526023	Keyword_counter	NIF	Information
20	3.162666096	medPause	IF	Unintentional Silence
21	3.145025579	postFillerPOS_Verb	IF	Filler
22	2.922614281	Verb_Ct_3sing	NIF	Semantic
23	2.770281656	postFillerPOS_Other	IF	Filler

Fig 5. Chart based on Appendix Table 6 data (Predictive Weight grouped by Factor/Umbrella)

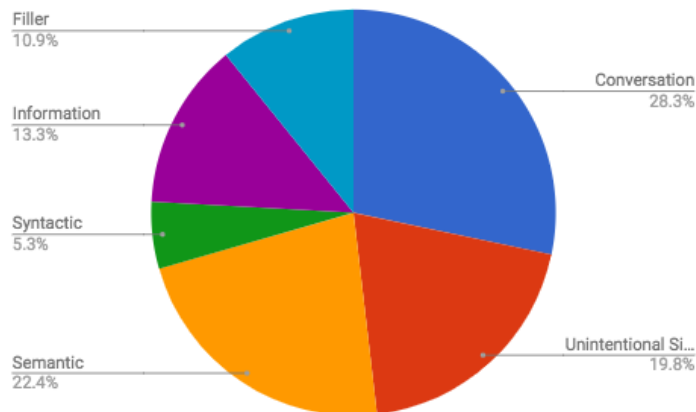


Fig 6. Chart based on Appendix Table 6 data (Count of final features grouped by Factor/Umbrella) **Ignoring weight value*

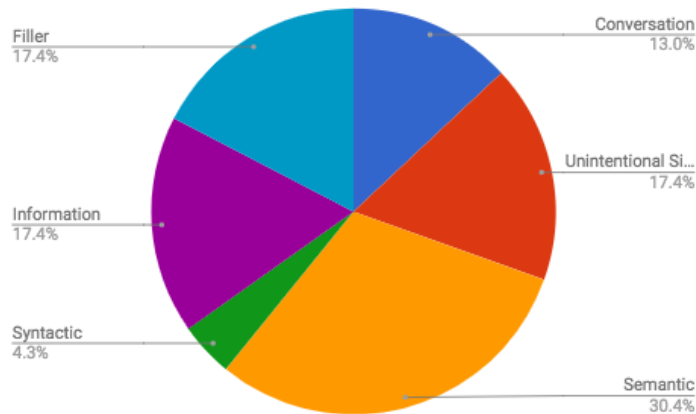


Table 7. Fraser et al. Features that this projects Non-IF's are based on.

Note. Below table is reprinted from [Fraser et al. 2016] (Table 2.)

Correlations with diagnosis (first column) and promax factor loadings.

Loadings less than 0.1 are excluded. Bold font indicates a loading greater than 0.3.

Feature	r	*Fac. 1	*Fac. 2	*Fac. 3	*Fac. 4
Pronoun:noun ratio	0.35	1.01		-0.32	
NP -> PRP	0.37	0.88		-0.24	
Frequency	0.34	0.74			
Adverbs	0.31	0.51			0.19
ADVP -> RB	0.3	0.44		0.1	
Verb frequency	0.21	0.39		0.13	
Nouns	-0.27	-0.97		0.37	
Word length	-0.41	-0.6		-0.13	
NP -> DT NN	0.1	-0.52			0.19
Honore's statistic	-0.25	-0.46	-0.14	-0.14	0.33
Inflected verbs	-0.19	-0.39			
Average cosine distance	-0.19	0.33		-0.15	0.13
Skewness(MFCC 1)	0.22		0.95		
Skewness(MFCC 2)	0.2		0.87		-0.14
Kurtosis(MFCC 5)	0.19		0.78		
Kurtosis(VEL(MFCC 3))	0.24	-0.17	0.44		0.24
Phonation rate	-0.21	0.16	-0.62		-0.28
Skewness(MFCC 8)	-0.22		-0.39		-0.13
Not-in-dictionary	0.38	-0.14		0.53	0.26
ROOT -> FRAG	0.23	-0.15		0.36	0.19
Verbs	-0.29	0.38		-1.05	0.2
VP rate	-0.19	0.37		-0.95	0.32
VP -> AUX VP	-0.23	0.16		-0.56	0.18
VP -> VBG	-0.27	-0.28		-0.34	0.21
Key word: window	-0.29			0.2	-0.79
Info unit: window	-0.32			0.12	-0.63
KEY WORD: sink	-0.23				-0.62
KEY WORD: cookie	-0.23		0.13		-0.61
PP proportion	-0.21			0.18	-0.61
Key word: curtain	-0.25				-0.56
PP rate	-0.21			0.19	-0.55
Info unit: curtain	-0.26				-0.53
Key word: counter	-0.18			0.14	-0.47
Info unit: cookie	-0.24				-0.46
Info unit: sink	-0.31				-0.43
Info unit: girl	-0.30				-0.42
Info unit: girl's action	-0.25		0.13	-0.12	-0.36
Info unit: dish	-0.24	-0.12			-0.29
Key word: stool	-0.28	-0.15			-0.29
Key word: mother	-0.32			-0.27	-0.26
Info unit: stool	-0.32	-0.29			-0.21
Skewness(MFCC 12)	-0.19				-0.18
Info unit: woman	-0.29			-0.16	-0.18
VP -> VBG PP	-0.34	-0.19		-0.30	-0.12
VP->IN S	-0.20				-0.1
VP -> AUX ADJP	-0.19	-0.11			
VP -> AUX	0.2	0.28			
VP -> VBD NP	0.19				
Cosine cutoff: 0.5	0.19		0.15	0.14	
INTJ -> UH	0.18	0.25			

*Fac. 1: Semantic Impairment, Fac. 2: Acoustic Impairment, Fac. 3 Syntactic Impairment, Fac. 4: Information Content Impairment