

# Content-based Frailty Classification System ++

Hospital Project

Project type: R&D (Competitive Grant)

Week 12 Final Presentation

Done By: Lee Jie Ying, Lynn (161015H), Low Ren Ern (163027J)



# Overview of Presentation

- Project Scope
- Clinical Frailty Scale
- Handling 2018 Old and New Sample Data from TTSH
- Difference between Old and New Sample Data
- Timeline (Week 1-5)
- Timeline (Week 6-11)
- Task Accomplished
- Flowchart of Project
- Problems Encountered
- Solutions
- Summing-up



# Project Scope

- This project is awarded under NYP-TTSH grant (18 Months)
- The purpose of this project is to develop a system to categorize the patient's notes to different frailty score.
- This is to improve the consistency of the score assignment so that they can render the suitable healthcare to the patient with the category.
- TTSH have provided 2015 and 2018 sample data that have been categorized for us to train the data.
  - Focus mainly on 2018 sample data



# Clinical Frailty Scale (CFS)

- Objective: To examine the CFS in patients admitted to acute medical ward and its association with length of stay.
- Frailty Score (also known as category in this project)



# Clinical Frailty Scale (CFS)

## Clinical Frailty Scale\*

For example: Category 1 is extremely fit and Category 9 is terminally ill



**1 Very Fit** – People who are robust, active, energetic and motivated. These people commonly exercise regularly. They are among the fittest for their age.



**2 Well** – People who have **no active disease symptoms** but are less fit than category 1. Often, they exercise or are very **active occasionally**, e.g. seasonally.



**3 Managing Well** – People whose **medical problems are well controlled**, but are **not regularly active** beyond routine walking.



**4 Vulnerable** – While **not dependent** on others for daily help, often **symptoms limit activities**. A common complaint is being “slowed up”, and/or being tired during the day.



**5 Mildly Frail** – These people often have **more evident slowing**, and need help in **high order IADLs** (finances, transportation, heavy housework, medications). Typically, mild frailty progressively impairs shopping and walking outside alone, meal preparation and housework.



**6 Moderately Frail** – People need help with **all outside activities** and with **keeping house**. Inside, they often have problems with stairs and need **help with bathing** and might need minimal assistance (cuing, standby) with dressing.



**7 Severely Frail** – **Completely dependent for personal care**, from whatever cause (physical or cognitive). Even so, they seem stable and not at high risk of dying (within ~ 6 months).



**8 Very Severely Frail** – Completely dependent, approaching the end of life. Typically, they could not recover even from a minor illness.



**9 Terminally Ill** - Approaching the end of life. This category applies to people with a **life expectancy <6 months**, who are **not otherwise evidently frail**.

### Scoring frailty in people with dementia

The degree of frailty corresponds to the degree of dementia. Common **symptoms in mild dementia** include forgetting the details of a recent event, though still remembering the event itself, repeating the same question/story and social withdrawal.

In **moderate dementia**, recent memory is very impaired, even though they seemingly can remember their past life events well. They can do personal care with prompting.

In **severe dementia**, they cannot do personal care without help.

\* 1. Canadian Study on Health & Aging, Revised 2008.

2. K. Rockwood et al. A global clinical measure of fitness and frailty in elderly people. CMAJ 2005;173:489-495.



# Handling 2018 Old Sample Data from TTSH

## ED (Emergency) Notes Sample Data

Total: **44** files

860246760\_Cat\_5\_ED Notes  
861236598\_Cat\_4\_ED Notes  
862832645\_Cat\_6\_ED Notes  
866043327\_Cat\_7\_ED Notes  
866972291\_Cat\_6\_ED Notes  
868528977\_Cat\_5\_ED Notes  
869240003\_Cat\_7\_ED Notes  
872877037\_Cat\_6\_ED Notes  
883716700\_Cat\_6\_ED Notes

860305644\_Cat\_5\_ED Notes  
861341492\_Cat\_5\_ED Notes  
864605611\_Cat\_6\_ED Notes  
866197419\_Cat\_3-4\_ED Notes  
867167556\_Cat\_4\_ED Notes  
868533852\_Cat\_5\_ED Notes  
869643579\_Cat\_7\_ED Notes  
882668111\_Cat\_8\_ED Notes  
883918817\_Cat\_4\_ED Notes

860389965\_Cat\_5\_ED Notes  
861541971\_Cat\_5\_ED Notes  
864704750\_Cat\_6\_ED Notes  
866312687\_Cat\_3\_ED Notes  
867359638\_Cat\_6\_ED Notes  
868589632\_Cat\_6-7\_ED Notes  
869797158\_Cat\_8\_ED Notes  
882695046\_Cat\_7\_ED Notes  
889029982\_Cat\_6\_ED Notes

860512095\_Cat\_5\_ED Notes  
861625560\_Cat\_7\_ED Notes  
864773682\_Cat\_7\_ED Notes  
866460588\_Cat\_6\_ED Notes  
867808171\_Cat\_4\_ED Notes  
868981160\_Cat\_6\_ED Notes  
872566043\_Cat\_5\_ED Notes  
882754748\_Cat\_4\_ED Notes  
892486974\_Cat\_7\_ED Notes

860755305\_Cat\_3\_ED Notes  
861864723\_Cat\_6\_ED Notes  
864873792\_Cat\_7\_ED Notes  
866745356\_Cat\_7\_ED Notes  
868194296\_Cat\_5\_ED Notes  
869095984\_Cat\_6\_ED Notes  
872871395\_Cat\_5\_ED Notes  
883054456\_Cat\_7\_ED Notes

## IFA Sample Data

Total: **48** files

860246760\_Cat\_5\_IFA  
861236598\_Cat\_4\_IFA  
862832645\_Cat\_6\_IFA  
866197419\_Cat\_3-4\_IFA  
867167556\_Cat\_4\_IFA  
868528977\_Cat\_5\_IFA  
869095984\_Cat\_6\_IFA  
872871395\_Cat\_5\_IFA  
882695046\_Cat\_7\_IFA  
889029982\_Cat\_6\_IFA

860305644\_Cat\_5\_IFA  
861341492\_Cat\_5\_IFA  
864605611\_Cat\_6\_IFA  
866312687\_Cat\_3\_IFA  
867359638\_Cat\_6\_IFA  
868533852\_Cat\_5\_IFA  
869240003\_Cat\_7\_IFA  
872877037\_Cat\_6\_IFA  
882754748\_Cat\_4\_IFA  
892486974\_Cat\_7\_IFA

860389965\_Cat\_5\_IFA  
861541971\_Cat\_5\_IFA  
864704750\_Cat\_6\_IFA  
866460588\_Cat\_6\_IFA  
867773642\_Cat\_7\_IFA  
868589632\_Cat\_6-7\_IFA  
869643579\_Cat\_7\_IFA  
875352977\_Cat\_4-5\_IFA  
883054456\_Cat\_7\_IFA

860512095\_Cat\_5\_IFA  
861625560\_Cat\_7\_IFA  
864773682\_Cat\_7\_IFA  
866745356\_Cat\_7\_IFA  
867808171\_Cat\_4\_IFA  
868981160\_Cat\_6\_IFA  
869797158\_Cat\_8\_IFA  
87996159\_Cat\_5-6\_IFA  
883716700\_Cat\_6\_IFA

860755305\_Cat\_3\_IFA  
861864723\_Cat\_6\_IFA  
864873792\_Cat\_7\_IFA  
866972291\_Cat\_6\_IFA  
868194296\_Cat\_5\_IFA  
869060647\_Cat\_4\_IFA  
872566043\_Cat\_5\_IFA  
882668111\_Cat\_8\_IFA  
883918817\_Cat\_4\_IFA





# Handling 2018 Old Sample Data from TTSH

Patient Care Record Sample Data (consists of Admission & Discharge)

Total: 82 files

860305644_Cat_5_Patient Care Record (Inpatien...	860305644_Cat_5_Patient Care Record (Inpatien...	860389965_Cat_5_Patient Care Record (Inpatien...	860512095_Cat_5_Patient Care Record (Inpatien...	860755305_Cat_3_Patient Care Record (Inpatien...
860755305_Cat_3_Patient Care Record (Inpatien...	861236598_Cat_4_Patient Care Record (Inpatien...	861236598_Cat_4_Patient Care Record (Inpatien...	861341492_Cat_5_Patient Care Record (Inpatien...	861541971_Cat_5_Patient Care Record (Inpatien...
861541971_Cat_5_Patient Care Record (Inpatien...	861625560_Cat_7_Patient Care Record (Inpatien...	861625560_Cat_7_Patient Care Record (Inpatien...	861864723_Cat_6_Patient Care Record (Inpatien...	861864723_Cat_6_Patient Care Record (Inpatien...
862832645_Cat_6_Patient Care Record (Inpatien...	864605611_Cat_6_Patient Care Record (Inpatie...	864605611_Cat_6_Patient Care Record (Inpatie...	864704750_Cat_6_Patient Care Record (Inpatien...	864704750_Cat_6_Patient Care Record (Inpatien...
864773682_Cat_7_Patient Care Record (Inpatien...	864773682_Cat_7_Patient Care Record (Inpatien...	864873792_Cat_7_Patient Care Record (Inpatien...	866043327_Cat_7_Patient Care Record (Inpatien...	866043327_Cat_7_Patient Care Record (Inpatien...
866197419_Cat_3-4_Patient Care Record (Inpat...	866197419_Cat_3-4_Patient Care Record (Inpat...	866312687_Cat_3_Patient Care Record (Inpatien...	866312687_Cat_3_Patient Care Record (Inpatien...	866460588_Cat_6_Patient Care Record (Inpatien...
866460588_Cat_6_Patient Care Record (Inpatien...	866745356_Cat_7_Patient Care Record (Inpatien...	866745356_Cat_7_Patient Care Record (Inpatien...	866972291_Cat_6_Patient Care Record (Inpatien...	867167556_Cat_4_Patient Care Record (Inpatien...
867167556_Cat_4_Patient Care Record (Inpatien...	867773642_Cat_7_Patient Care Record (Inpatien...	867773642_Cat_7_Patient Care Record (Inpatien...	867808171_Cat_4_Patient Care Record (Inpatien...	867808171_Cat_4_Patient Care Record (Inpatien...
868194296_Cat_5_Patient Care Record (Inpatien...	868194296_Cat_5_Patient Care Record (Inpatien...	868528977_Cat_5_Patient Care Record (Inpatien...	868528977_Cat_5_Patient Care Record (Inpatien...	868533852_Cat_5_Patient Care Record (Inpatien...
868533852_Cat_5_Patient Care Record (Inpatien...	868589632_Cat_6-7_Patient Care Record (Inpat...	868589632_Cat_6-7_Patient Care Record (Inpat...	868981160_Cat_6_Patient Care Record (Inpatien...	868981160_Cat_6_Patient Care Record (Inpatien...
869060647_Cat_4_Patient Care Record (Inpatien...	869060647_Cat_4_Patient Care Record (Inpatien...	869095984_Cat_6_Patient Care Record (Inpatien...	869240003_Cat_7_Patient Care Record (Inpatien...	869240003_Cat_7_Patient Care Record (Inpatien...
869643579_Cat_7_Patient Care Record (Inpatien...	869643579_Cat_7_Patient Care Record (Inpatien...	869797158_Cat_8_Patient Care Record (Inpatien...	869797158_Cat_8_Patient Care Record (Inpatien...	872566043_Cat_5_Patient Care Record (Inpatien...
872566043_Cat_5_Patient Care Record (Inpatien...	872871395_Cat_5_Patient Care Record (Inpatien...	872871395_Cat_5_Patient Care Record (Inpatien...	872877037_Cat_6_Patient Care Record (Inpatien...	872877037_Cat_6_Patient Care Record (Inpatien...
875352977_Cat_4-5_Patient Care Record (Inpat...	875352977_Cat_4-5_Patient Care Record (Inpat...	879996159_Cat_5-6_Patient Care Record (Inpat...	879996159_Cat_5-6_Patient Care Record (Inpat...	882668111_Cat_8_Patient Care Record (Inpatien...
882668111_Cat_8_Patient Care Record (Inpatien...	882695046_Cat_7_Patient Care Record (Inpatien...	882754748_Cat_4_Patient Care Record (Inpatien...	882754748_Cat_4_Patient Care Record (Inpatien...	883054456_Cat_7_Patient Care Record (Inpatien...
883054456_Cat_7_Patient Care Record (Inpatien...	883716700_Cat_6_Patient Care Record (Inpatien...	883918817_Cat_4_Patient Care Record (Inpatien...	883918817_Cat_4_Patient Care Record (Inpatien...	889029982_Cat_6_Patient Care Record (Inpatien...
889029982_Cat_6_Patient Care Record (Inpatien...	892486974_Cat_7_Patient Care Record (Inpatien...			



# Handling 2018 New Sample Data from TTSH

ED (Emergency) Notes  
Sample Data

Total: **45** files

560277507\_Cat\_7\_ED Notes  
576867505\_Cat\_6\_ED Notes  
585496133\_Cat\_5\_ED Notes  
679020725\_Cat\_7\_ED Notes  
862166495\_Cat\_5\_ED Notes  
862323333\_Cat\_7\_ED Notes  
862349405\_Cat\_6\_ED Notes  
862456463\_Cat\_5\_ED Notes  
862649075\_Cat\_5\_ED Notes

561005678\_Cat\_6\_ED Notes  
578935973\_Cat\_7\_ED Notes  
663198575\_Cat\_7\_ED Notes  
862033106\_Cat\_3-4\_ED Notes  
862232907\_Cat\_6\_ED Notes  
862328287\_Cat\_6\_ED Notes  
862358364\_Cat\_3\_ED Notes  
862462779\_Cat\_5\_ED Notes  
862810017\_Cat\_4\_ED Notes

568831565\_Cat\_6\_ED Notes  
579345112\_Cat\_8\_ED Notes  
667598022\_Cat\_7\_ED Notes  
862073014\_Cat\_6\_ED Notes  
862254714\_Cat\_4\_ED Notes  
862342662\_Cat\_5\_ED Notes  
862367768\_Cat\_6\_ED Notes  
862465085\_Cat\_8\_ED Notes  
862861782\_Cat\_5\_ED Notes

571192553\_Cat\_2\_ED Notes  
579514165\_Cat\_5\_ED Notes  
668051254\_Cat\_5\_ED Notes  
862096578\_Cat\_4\_ED Notes  
862270808\_Cat\_7\_ED Notes  
862345159\_Cat\_8\_ED Notes  
862407657\_Cat\_6\_ED Notes  
862469994\_Cat\_6\_ED Notes  
862875722\_Cat\_3-4\_ED Notes

575911924\_Cat\_6\_ED Notes  
584121291\_Cat\_6\_ED Notes  
677153773\_Cat\_6\_ED Notes  
862096586\_Cat\_4\_ED Notes  
862298725\_Cat\_2\_ED Notes  
862348514\_Cat\_6\_ED Notes  
862408876\_Cat\_5\_ED Notes  
862572175\_Cat\_5\_ED Notes  
862895545\_Cat\_3\_ED Notes

IFA Sample Data

Total: **92** files ??





# Handling 2018 New Sample Data from TTSH

Patient Care Record Sample Data (consists of Admission & Discharge)

Total: 92 files

560277507_Cat_7_Patient Care Record (Inpatien...	560277507_Cat_7_Patient Care Record (Inpatien...	561005678_Cat_6_Patient Care Record (Inpatien...	561005678_Cat_6_Patient Care Record (Inpatien...	568831565_Cat_6_Patient Care Record (Inpatien...
568831565_Cat_6_Patient Care Record (Inpatien...	571192553_Cat_2_Patient Care Record (Inpatien...	571192553_Cat_2_Patient Care Record (Inpatien...	575911924_Cat_6_Patient Care Record (Inpatien...	575911924_Cat_6_Patient Care Record (Inpatien...
576867505_Cat_6_Patient Care Record (Inpatien...	576867505_Cat_6_Patient Care Record (Inpatien...	578935973_Cat_7_Patient Care Record (Inpatien...	578935973_Cat_7_Patient Care Record (Inpatien...	579345112_Cat_8_Patient Care Record (Inpatien...
579345112_Cat_8_Patient Care Record (Inpatien...	579514165_Cat_5_Patient Care Record (Inpatien...	579514165_Cat_5_Patient Care Record (Inpatien...	584121291_Cat_6_Patient Care Record (Inpatien...	584121291_Cat_6_Patient Care Record (Inpatien...
585496133_Cat_5_Patient Care Record (Inpatien...	585496133_Cat_5_Patient Care Record (Inpatien...	663198575_Cat_7_Patient Care Record (Inpatien...	663198575_Cat_7_Patient Care Record (Inpatien...	667598022_Cat_7_Patient Care Record (Inpatien...
667598022_Cat_7_Patient Care Record (Inpatien...	668051254_Cat_5_Patient Care Record (Inpatien...	668051254_Cat_5_Patient Care Record (Inpatien...	677153773_Cat_6_Patient Care Record (Inpatien...	677153773_Cat_6_Patient Care Record (Inpatien...
679020725_Cat_7_Patient Care Record (Inpatien...	679020725_Cat_7_Patient Care Record (Inpatien...	862033106_Cat_3-4_Patient Care Record (Inpati...	862033106_Cat_3-4_Patient Care Record (Inpati...	862073014_Cat_6_Patient Care Record (Inpatien...
862073014_Cat_6_Patient Care Record (Inpatien...	862096578_Cat_4_Patient Care Record (Inpatien...	862096578_Cat_4_Patient Care Record (Inpatien...	862096586_Cat_4_Patient Care Record (Inpatien...	862096586_Cat_4_Patient Care Record (Inpatien...
862166495_Cat_5_Patient Care Record (Inpatien...	862166495_Cat_5_Patient Care Record (Inpatien...	862232907_Cat_6_Patient Care Record (Inpatien...	862232907_Cat_6_Patient Care Record (Inpatien...	862254714_Cat_4_Patient Care Record (Inpatien...
862254714_Cat_4_Patient Care Record (Inpatien...	862270808_Cat_7_Patient Care Record (Inpatien...	862270808_Cat_7_Patient Care Record (Inpatien...	862298725_Cat_2_Patient Care Record (Inpatien...	862298725_Cat_2_Patient Care Record (Inpatien...
862323333_Cat_7_Patient Care Record (Inpatien...	862323333_Cat_7_Patient Care Record (Inpatien...	862328287_Cat_6_Patient Care Record (Inpatien...	862328287_Cat_6_Patient Care Record (Inpatien...	862342662_Cat_5_Patient Care Record (Inpatien...
862342662_Cat_5_Patient Care Record (Inpatien...	862345159_Cat_8_Patient Care Record (Inpatien...	862345159_Cat_8_Patient Care Record (Inpatien...	862348514_Cat_6_Patient Care Record (Inpatien...	862348514_Cat_6_Patient Care Record (Inpatien...
862349405_Cat_6_Patient Care Record (Inpatien...	862349405_Cat_6_Patient Care Record (Inpatien...	862358364_Cat_3_Patient Care Record (Inpatien...	862358364_Cat_3_Patient Care Record (Inpatien...	862367768_Cat_6_Patient Care Record (Inpatien...
862367768_Cat_6_Patient Care Record (Inpatien...	862407657_Cat_6_Patient Care Record (Inpatien...	862407657_Cat_6_Patient Care Record (Inpatien...	862408876_Cat_5_Patient Care Record (Inpatien...	862408876_Cat_5_Patient Care Record (Inpatien...
862456463_Cat_5_Patient Care Record (Inpatien...	862456463_Cat_5_Patient Care Record (Inpatien...	862462779_Cat_5_Patient Care Record (Inpatien...	862462779_Cat_5_Patient Care Record (Inpatien...	862465085_Cat_8_Patient Care Record (Inpatien...
862465085_Cat_8_Patient Care Record (Inpatien...	862469994_Cat_6_Patient Care Record (Inpatien...	862469994_Cat_6_Patient Care Record (Inpatien...	862572175_Cat_5_Patient Care Record (Inpatien...	862572175_Cat_5_Patient Care Record (Inpatien...
862649075_Cat_5_Patient Care Record (Inpatien...	862649075_Cat_5_Patient Care Record (Inpatien...	862810017_Cat_4_Patient Care Record (Inpatien...	862810017_Cat_4_Patient Care Record (Inpatien...	862861782_Cat_5_Patient Care Record (Inpatien...
862861782_Cat_5_Patient Care Record (Inpatien...	862875722_Cat_3-4_Patient Care Record (Inpati...	862875722_Cat_3-4_Patient Care Record (Inpati...	862894859_Cat_6_Patient Care Record (Inpatien...	862894859_Cat_6_Patient Care Record (Inpatien...
862895545_Cat_3_Patient Care Record (Inpatien...	862895545_Cat_3_Patient Care Record (Inpatien...			



# Difference between Old and New Sample Data

## Old Sample Data

- Documents have been converted before we handle them
- Word documents' data format is inconsistent
- Example of inconsistency in data format:
  - Page breaks
  - Missing values in tables

## New Sample Data

- Documents were in PDF and not converted
- Word documents have to be converted to PDF files
- Newer sample data format is more consistent

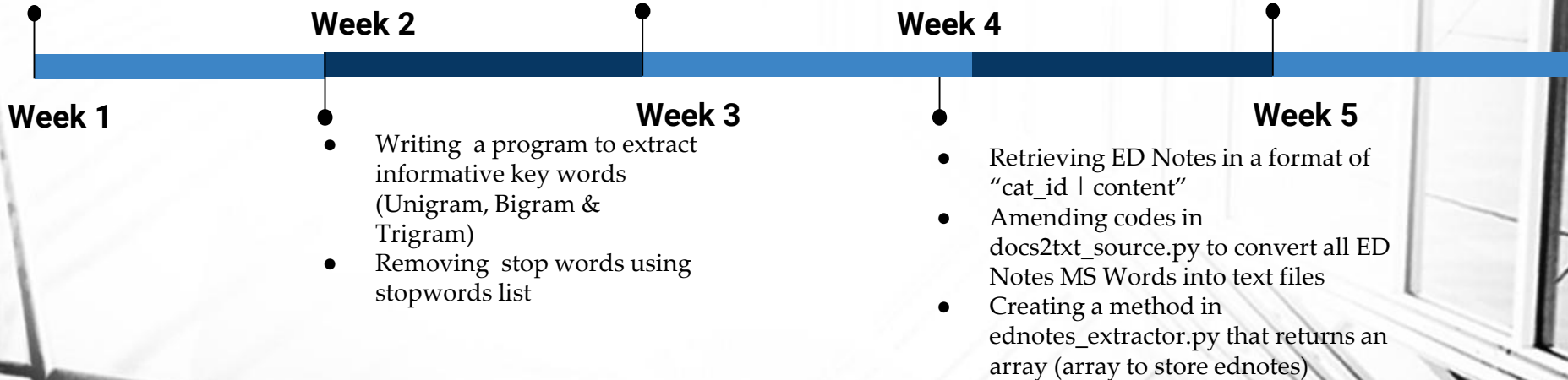


## Timeline (Week 1-5) - Lynn

- Learning of Python & MongoDB
- Retrieving ED Notes from MongoDB
- Understanding text mining

- Retrieving of ED Notes to store in an array for text analysis later
- Improving Unigram, Bigram & Trigram

- Trying to finish up docs2txt\_source.py for ED Notes
- Combine my two testing projects together into one project





## Timeline (Week 6-11) - Lynn

- Create a program that handles 3 different input files (ED Notes, Patient Care Records Admission, Patient Care Records Discharge)
  - Create inpatientcare extractor to extract Patient Care Records Admission and Discharge
  - Trying another approach to improve the mean accuracy of data sets
  - Learning of Topic Modeling and trying the practical for topic modeling
  - Extract word terms from LDA (Latent Dirichlet Allocation) model in topic modeling
  - Store extracted word terms in a dictionary and append to ngram key phrases list
- 
- The timeline is represented by a horizontal bar divided into segments for each week. Markers are placed at the start of each week segment, with task descriptions listed to the left or right of the bar.
- | Week     | Tasks   |
|----------|---|
| Week 6-7 | <ul style="list-style-type: none"><li>• Create a program that handles 3 different input files (ED Notes, Patient Care Records Admission, Patient Care Records Discharge)</li><li>• Create inpatientcare extractor to extract Patient Care Records Admission and Discharge</li></ul> |
| Week 8   | <ul style="list-style-type: none"><li>• N-grams (Knowledge-based)<ul style="list-style-type: none"><li>- Using this for keyword matching</li></ul></li><li>• Handling new and old data sets for oversampling to improve the mean accuracy of data sets</li></ul>                    |
| Week 9   | <ul style="list-style-type: none"><li>• Trying another approach to improve the mean accuracy of data sets</li><li>• Learning of Topic Modeling and trying the practical for topic modeling</li></ul>  |
| Week 10  | <ul style="list-style-type: none"><li>• Trying to implement topic modeling into the current working project</li><li>• Create keywords_consensus_analysis.py to do keyword matching</li></ul>  |
| Week 11  | <ul style="list-style-type: none"><li>• Extract word terms from LDA (Latent Dirichlet Allocation) model in topic modeling</li><li>• Store extracted word terms in a dictionary and append to ngram key phrases list</li></ul>   |



## Tasks Accomplished - Lynn

- Create `inpatientcare_extractor.py` to extract Patient Care Record Discharge and Admission data directly from word documents.
- Create a program to handle 3 different input files (ED Notes, Patient Care Record\_Admission & Patient Care Record\_Discharge) for text analysis.
- Improve mean accuracy results by using oversampled data (Oversampling of imbalanced data)
- Create `key_consensus_analysis.py` to do keyword matching
- Create `topicmodel_keywords_extractor.py` to extract keywords from LDA (Latent Dirichlet Allocation) model in topic modeling

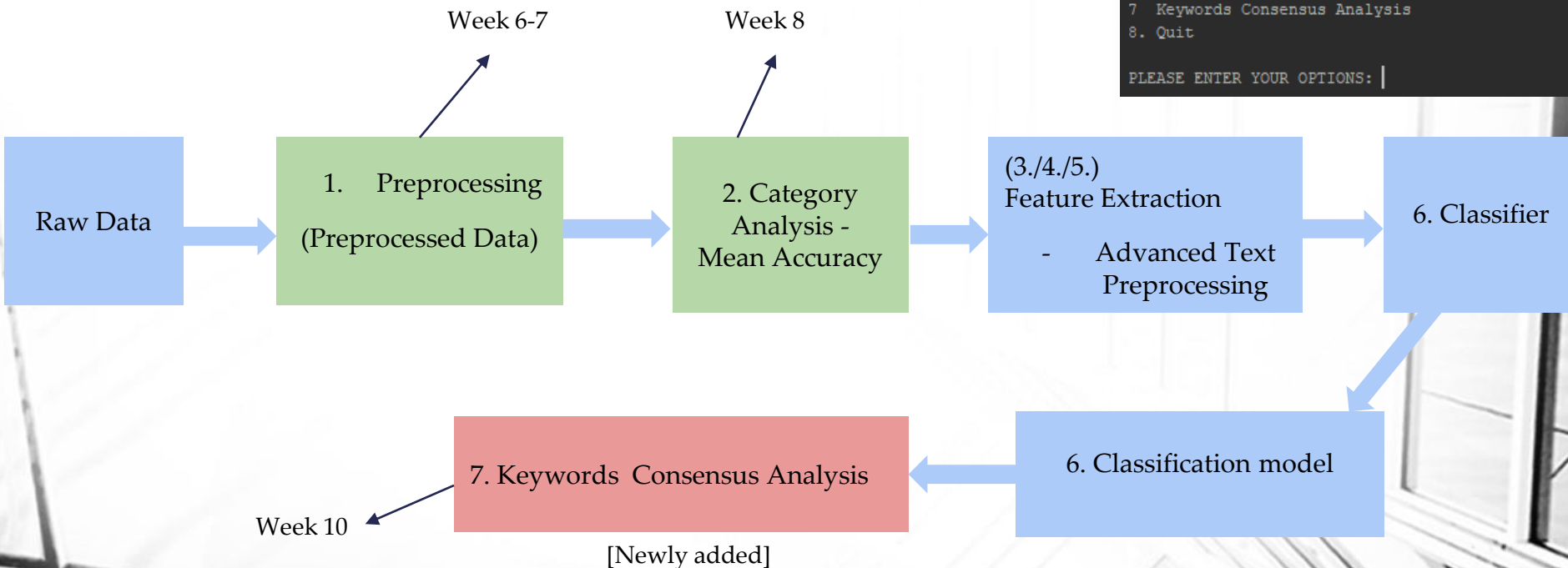


# Flowchart of Project - Lynn

```
===== MENU =====
=====

OPTIONS:
1. Pre-processing - DocsToTxt
2. Category Analysis - Mean Accuracy
3. Unigram
4. Unigram & Bigram
5. Unigram, Bigram & Trigram
6. Co-efficient for 4 models
7. Keywords Consensus Analysis
8. Quit

PLEASE ENTER YOUR OPTIONS: |
```







# Flowchart of Project - Lynn

Raw Data

ED Notes, Patient Care  
Record Admission, Patient  
Care Record Discharge

1. Preprocessing  
(Preprocessed Data)

Week 6-7

Basic preprocessing of text data

Go to docs2txt\_source.py that  
handles all 3 different output files  
(ED Notes, Patient Care Record  
Admission & Patient Care Record  
Discharge)

Yes

Remove common  
phases/headers

No

Method -  
get\_ednotes(source\_directory) in  
ednotes\_extractor.py called by  
docs2txt\_source.py

Goes to source directory, gets each  
ed notes file, stores in an array & it  
returns an array

Save all output files into  
docs2txt\_output.txt file  
in "cat\_id|content"  
format

```
===== MENU =====
=====
OPTIONS:
1. Preprocessing - DocsToTxt
2. CategoryAnalysis - Mean Accuracy
3. Unigram
4. Unigram & Bigram
5. Unigram, Bigram & Trigram
6. Co-efficient for 4 models
7. Quit
```

```
PLEASE ENTER YOUR OPTIONS: 1
```

```
REMOVING COMMON PHASES, HEADERS? Y/N: Y
```

```
▼ dataprep
  ► documents_sampling
  ► source_documents
  ► docs2txt_output.txt
```



# Flowchart of Project - Lynn

1. Preprocessing  
(Preprocessed Data)

2. Category Analysis -  
Mean Accuracy

Week 8

## Category Analysis - Mean Accuracy

The method - `mean_acc()` in the `category_analysis.py` is being called

```
model_name
LinearSVC      0.340540
LogisticRegression  0.351332
MultinomialNB  0.330268
RandomForestClassifier 0.298286
Name: accuracy, dtype: float64
```

\*\*\*Note: Oversampling for old Patient Care Records Admission and Discharge only

Before Oversampling:

```
model_name
LinearSVC      0.350736
LogisticRegression  0.360346
MultinomialNB  0.351688
RandomForestClassifier 0.399394
```

After Oversampling:

```
model_name
LinearSVC      0.757010
LogisticRegression  0.751377
MultinomialNB  0.615189
RandomForestClassifier 0.690137
```

```
===== MENU =====
=====
OPTIONS:
1. Preprocessing - DocsToTxt
2. CategoryAnalysis - Mean Accuracy
3. Unigram
4. Unigram & Bigram
5. Unigram, Bigram & Trigram
6. Co-efficient for 4 models
7. Quit

PLEASE ENTER YOUR OPTIONS: 2
```



# Flowchart of Project - Lynn

2. Category Analysis - Mean Accuracy



3./4./5. Feature Extraction

- Advanced Text Preprocessing

```
'''ngram_range(1, 1) -> unigram / ngram_range(1, 3) -> unigram, bigram, trigram'''  
tfidf = TfidfVectorizer(sublinear_tf=True, min_df=1, norm='l2', encoding='UTF-8', ngram_range=(1, 3),  
                        stop_words=stop_words)
```

N-grams:

- Unigram - extract a single word
- Bigram - extract a word pair
- Trigram - extract triple words

For example: `ngram_range(1, 1)` means unigram,  
`ngram_range(1, 3)` means unigram, bigram and trigram

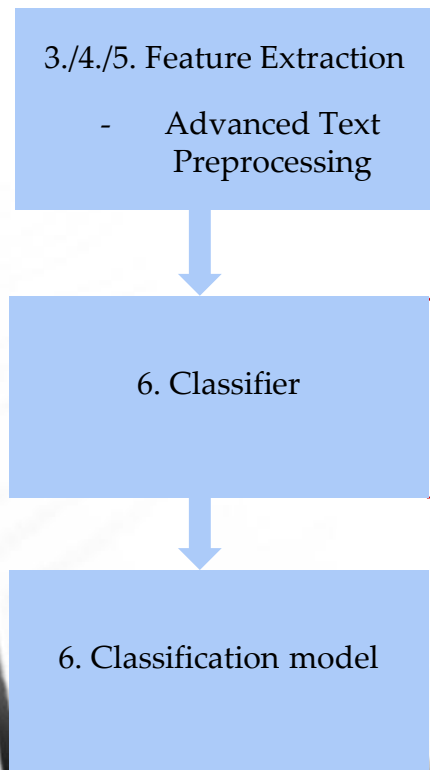
OPTIONS:

1. Preprocessing - DocsToTxt
2. CategoryAnalysis - Mean Accuracy
3. Unigram
4. Unigram & Bigram
5. Unigram, Bigram & Trigram
6. Co-efficient for 4 models
7. Quit

PLEASE ENTER YOUR OPTIONS: 5



# Flowchart of Project - Lynn



## 4 Classifiers:

- Naive Bayes (Multinomial)
- Logistic Regression
- Linear SVC
- Random Forest

RUNNING CATEGORY ANALYSIS...

SELECT A MODEL TYPE:

1. Naive Bayes
2. Logistic Regression
3. Linear SVC
4. Random Forest

PLEASE ENTER YOUR OPTION: |

```
=====
===== MENU =====
=====
OPTIONS:
1. Preprocessing - DocsToTxt
2. CategoryAnalysis - Mean Accuracy
3. Unigram
4. Unigram & Bigram
5. Unigram, Bigram & Trigram
6. Co-efficient for 4 models
7. Quit

PLEASE ENTER YOUR OPTIONS: 6
```

Classifier	Classification Model
An algorithm that maps input data to a specific category	Draws some conclusion from input values given for training



# Flowchart of Project - Lynn

Objective: Predict category from an unknown document string

## 7. Keywords Consensus Analysis

```
===== MENU =====
OPTIONS:
1. Pre-processing - DocsToTxt
2. Category Analysis - Mean Accuracy
3. Unigram
4. Unigram & Bigram
5. Unigram, Bigram & Trigram
6. Co-efficient for 4 models
7. Keywords Consensus Analysis
8. Quit

PLEASE ENTER YOUR OPTIONS: 7
```

5. Return the highest float value and the predicted category

Method - `check_keyword()` in `keywords_consensus_analysis.py`

1. Read an unknown document in a text file as a string

4. If keywords comparison found and match, then add count

The total count is then divided by the total number of keywords in the dictionary.

3. Compare unknown document string with the ngram keywords dictionary

2. Read `key_phrases` text file, which consists of ngram keywords as a dictionary



# Flowchart of Project - Lynn

Objective: Predict category from an unknown document string

## 7. Keywords Consensus Analysis

```
===== MENU =====
=====
OPTIONS:
1. Pre-processing - DocsToTxt
2. Category Analysis - Mean Accuracy
3. Unigram
4. Unigram & Bigram
5. Unigram, Bigram & Trigram
6. Co-efficient for 4 models
7. Keywords Consensus Analysis
8. Quit

PLEASE ENTER YOUR OPTIONS: ?
```

```
def check_keyword():
```

```
    # read unknown text from text file as a string
    with open(unknown, "r") as fp:
        unknown_file = fp.read()
```

```
    # read key phases from text file as a dictionary
    df = pd.read_csv(key_phrases, sep='|')
    key_phrases_dict = df.to_dict(orient='records')
```

```
    count = max(count_dict.values())
    # get the maximum value in the dictionary
    max_value = [(k, v) for k, v in count_dict.items() if v == count]
    maxval = print(max_value)
```

```
    # compare content of unknown file with key phases (write a for loop to do keyword matching)
    for key in new_dict.keys():
        count_dict[key] = 0
        new_list2 = new_dict[key].split(",")
        new_dict[key] = new_list2
        for j in new_dict[key]:
            if j in unknown_file:
                # if word found and matched, add count |
                count_dict[key] = count_dict[key] + 1
        count_dict[key] = float(count_dict[key] / len(new_list2))
```



## 7. Keywords Consensus Analysis

Objective: Predict category from an unknown document string

```
# compare content of unknown file with key phases (write a for loop to do keyword matching)
for key in new_dict.keys():
    count_dict[key] = 0
    new_list2 = new_dict[key].split(",")
    new_dict[key] = new_list2
    for j in new_dict[key]:
        if j in unknown_file:
            # if word found and matched, add count
            count_dict[key] = count_dict[key] + 1
    count_dict[key] = float(count_dict[key] / len(new_list2))
print(count_dict)
```

```
{2: 0.02666666666666667, 3: 0.006666666666666667, 4: 0.013333333333333334, 5: 0.04, 6: 0.08, 7: 0.03333333333333333, 8: 0.013333333333333334}
```

```
count = max(count_dict.values())
# get the maximum value in the dictionary
max_value = [(k, v) for k, v in count_dict.items() if v == count]
maxval = print(max_value)
```

[(6, 0.08)]

Predicted document to be category 6

Unknown document string is category 5

▼ converted\_documents  
579514165 Cat\_5\_Patient Care Record (Inpatient Nursin



# Flowchart of Project - Lynn

Objective: To add more informative words into key\_phrases.txt to improve the prediction of an unknown document string

## 7. Keywords Consensus Analysis

```
===== MENU =====
=====
OPTIONS:
1. Pre-processing - DocsToTxt
2. Category Analysis - Mean Accuracy
3. Unigram
4. Unigram & Bigram
5. Unigram, Bigram & Trigram
6. Co-efficient for 4 models
7. Keywords Consensus Analysis
8. Quit

PLEASE ENTER YOUR OPTIONS: 7
```

Method - get\_topicmodel\_words (filename) in topicmodel\_keywords\_extractor.py

Read in each category text file as a filename

```
get_topicmodel_words(output_cat3)
get_topicmodel_words(output_cat4)
get_topicmodel_words(output_cat5)
get_topicmodel_words(output_cat6)
get_topicmodel_words(output_cat7)
get_topicmodel_words(output_cat8)
```

key\_phrases to match cat\_id with the extracted word terms

If match then append to the words to the string of the particular category

Extract keywords from each category from LDA Model in topic modeling



# Problems Encountered - Lynn

- To handle new data sets given to us
- To handle old data sets with inconsistency in data format
- To improve on the Ngram's keyword extraction results



# Solutions - Lynn

- Seek help from previous batch students and supervisor
- Research online for solutions (e.g. stackoverflow)
- Trial and error by using several different methods
  - Different methods meaning to try out and test if online related codes to the project works

# Summing-up - Lynn

## **docs2txt\_source.py**

(for handling 3 different output files - ED Notes, Patient Care Record Admission & Patient Care Record Discharge)

To save all output files into docs2txt\_output.txt file in "cat\_id|content" format

## **ednotes\_extractor.py**

(for extracting ED Notes directly from the word documents)

To extract all ED Notes data text from the word documents

## **inpatientcare\_extractor**

(for extracting Patient Care Records Admission & Discharge directly from the word documents)

Extract method using docx2txt.process( )

import docx2txt: to extract text from docx files

## **keywords\_consensus\_analysis.py**

(for doing keyword matching between an unknown document string and a known dictionary of words)

To predict category from unknown document string

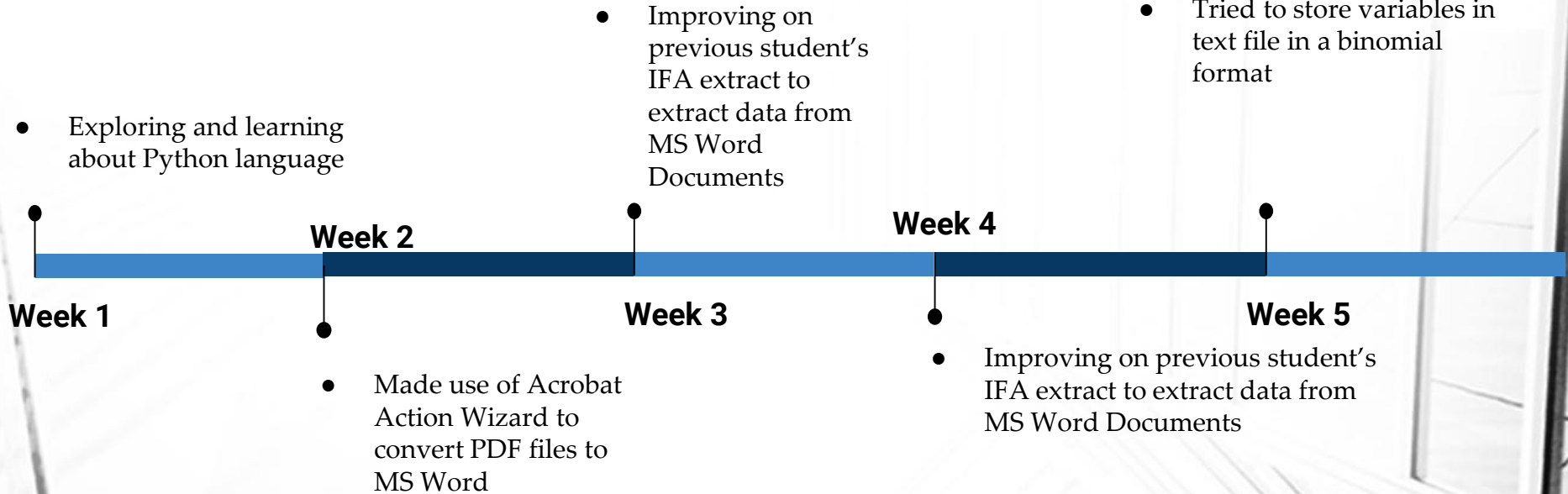
## **topicmodel\_keywords\_extractor.py**

(for extracting keywords for each category from LDA Model in Topic Modeling)

To add on keywords to the key\_phrases text file, which consists of ngram key phrases for each category



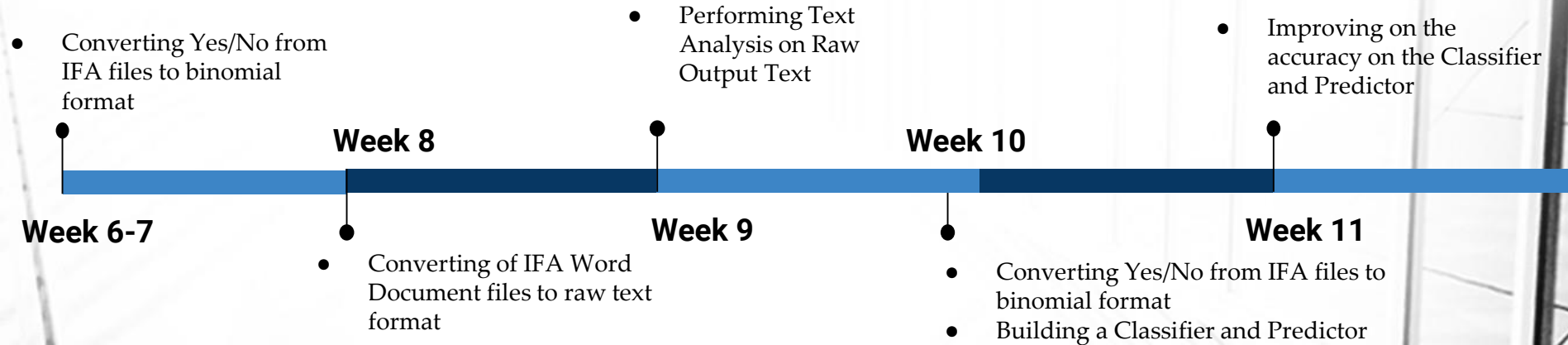
## Timeline (Week 1-5) - Ren Ern







## Timeline (Week 6-11) - Ren Ern





## Tasks Accomplished - Ren Ern

- Converting of IFA documents from PDF to Word Document using Adobe Acrobat
- Improving on the previous batch student's IFA extract to extract more data
- Converting IFA documents from Word Document format to Raw Output Text format
- Converting Yes/No from IFA files and converting them to binomial format
- Create a Classifier to identify relationships and patterns among the categories
- Create a Predictor to predict a category base on the Classifier



# Introduction

- IFA files contain medical information about a patient
- They are split into 5 categories
  1. Sensory and Communications
  2. Functional Status
  3. Mental Health
  4. Other Geriatric Syndromes
  5. Social

## DEPARTMENT OF GERIATRIC MEDICINE Inpatient Frailty Assessment (IFA)

### CHIEF COMPLAINT:

#### SENSORY and COMMUNICATION

Visual Impairment	Yes	Both
Hearing Impairment	Yes	Both
Communication	Yes	Verbal

#### FUNCTIONAL STATUS

**Decline in function  
compared to premorbid**

No

#### MENTAL HEALTH

Progressive Forgetfulness or Known Cognitive impairment	Yes	Advanced Dementia
Depression	Nil	
Sleep Issues	Yes	Sleep- wake reversal
Behavioral Disturbance	Nil	

Cognitive - known Advanced Dementia



# Introduction

- IFA files contain medical information about a patient
- They are split into 5 categories
  1. Sensory and Communications
  2. Functional Status
  3. Mental Health
  4. Other Geriatric Syndromes
  5. Social

## OTHER GERIATRIC SYNDROMES

*\*Kindly elaborate further if answer is YES*

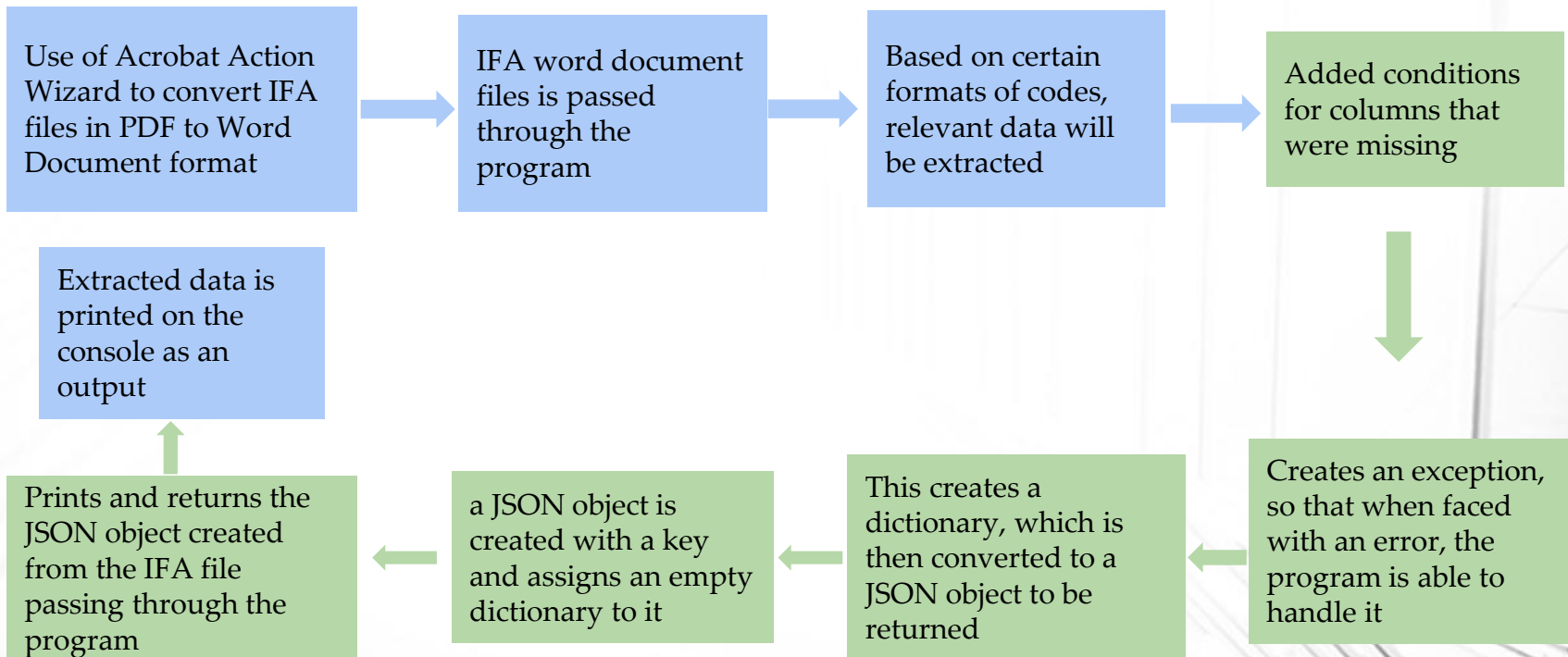
Swallowing Impairment	No	DOC normal texture with thin fluids
Loss of Appetite	No	
Unintentional Loss of Weight	No	
Urinary Retention	No	
Altered Bowel Habits	No	
Falls	Yes	Recurrent 2 falls this year, 1st on 1/5/18 fell on the grass on the way back from grocery shopping, claimed tripped and fell, was helped by <u>neighbours</u> and able to walk back to own home. 2nd fall 15/5/18 unwitnessed, <u>loss balance</u> while buying newspaper since second fall, noted functional decline and patient became <u>more weaker</u> . ADLs assisted and became homebound, required WC ambulation with assistance to transfer
Chronic Pain	No	

## SOCIAL

Smoking / Alcohol History	Nil	
Education Level	-	
Social History	<u>Orangevalley</u> nursing home resident Main <u>spokeperson</u> : Granddaughter <u>Ms Doris Mak</u>	
Caregiver Stress	No	



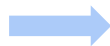
## Flowchart of Project - Ren Ern



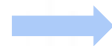


## Flowchart of Project - Ren Ern

Extracted data is printed on the console as an output



Create a method inside the program to convert the output of the console to a text file



Process each file



Store into an output



Output is dumped into a text file

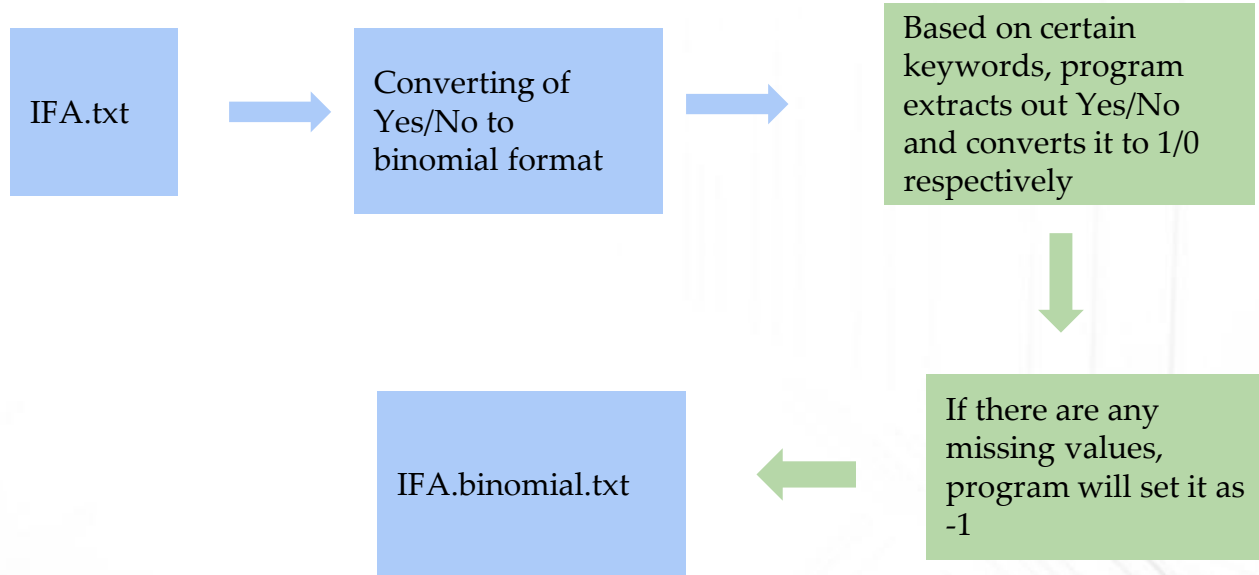


IFA.txt



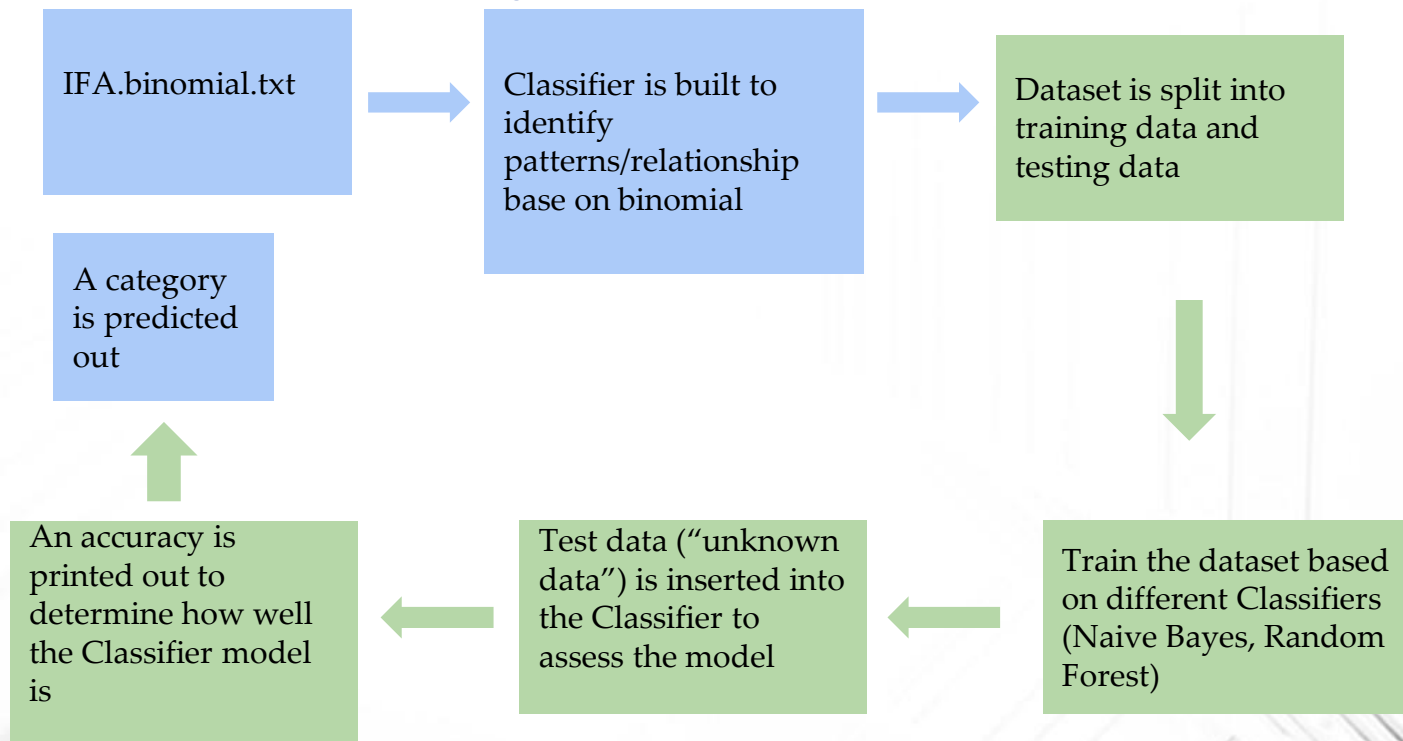


## Flowchart of Project - Ren Ern





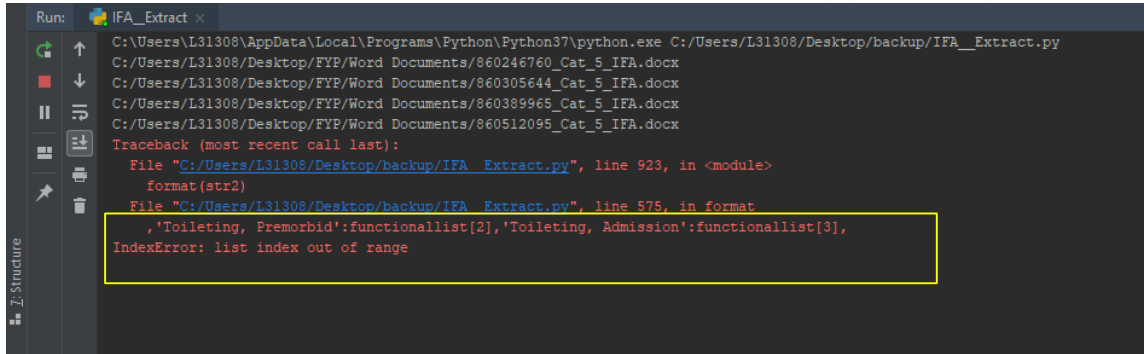
## Flowchart of Project - Ren Ern





# Problems Encountered - Ren Ern

- Previous student's extract was hard-coded to a single document, it was not “flexible” in reading different format of IFA documents



```
Run: IFA_Extract x
C:\Users\L31308\AppData\Local\Programs\Python\Python37\python.exe C:/Users/L31308/Desktop/backup/IFA_Extract.py
C:/Users/L31308/Desktop/FYP/Word Documents/860246760_Cat_5_IFA.docx
C:/Users/L31308/Desktop/FYP/Word Documents/860305644_Cat_5_IFA.docx
C:/Users/L31308/Desktop/FYP/Word Documents/860389965_Cat_5_IFA.docx
C:/Users/L31308/Desktop/FYP/Word Documents/860512095_Cat_5_IFA.docx
Traceback (most recent call last):
  File "C:/Users/L31308/Desktop/backup/IFA_Extract.py", line 923, in <module>
    format(str2)
  File "C:/Users/L31308/Desktop/backup/IFA_Extract.py", line 575, in format
    , 'Toileting, Premorbid':functionallist[2], 'Toileting, Admission':functionallist[3],
IndexError: list index out of range
```

Only 4 files were able to be read.



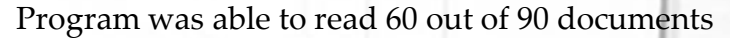
Previous student's extract was unable to read this document due to the page break

FUNCTIONAL STATUS	
Decline in function compared to premorbid	<p>No / Yes            (If YES indicate <i>function</i> on admission in the following functional assessment)            (following functional assessment)            assessment)</p> <p>Duration of decline if <i>ppts</i> if present:</p>
<p>*Premorbid = TWO weeks prior t prior to onset of acute illness or best function in last SIX months or best function in last SIX months months            Not applicable (e.g. <i>unconscious</i>, uncommunicative)</p>	
<p>I = Independent, A = Assisted/Supervised, Assisted/Supervised, D = Dependent  <u>Dependent</u></p>	

Created by: HERNANDEZ HERB HOWARD CUNANAN on 30-May-2018 18:39 on 9070  
 Printed by: AVRININI KRISHNA KISHORE on 30-Jun-2018 13:26

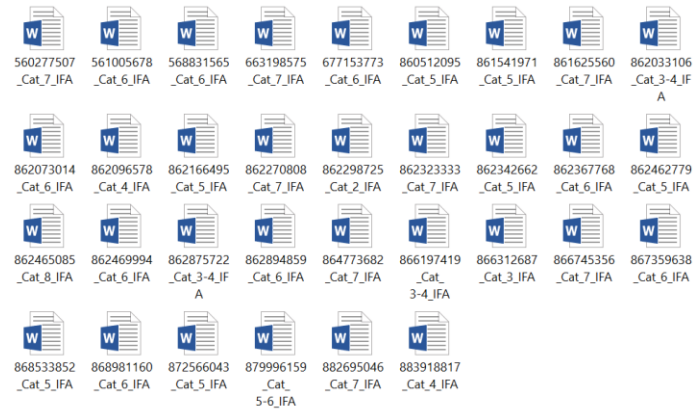
Page 2 of 7

*Kindly elaborate further if answer <u>stronger</u> is YES		
Swallowing Impairment	No / Yes	NG tube / PEG tube / tube / Modified consistency: consistency:
Loss of Appetite	No / Yes	





# Problems Encountered- Ren Ern



Documents that were still unable to be read



## Problems Encountered- Ren Ern

- Half of the Word Documents that were unable to be extracted had tables that were “cut”
- Hence program was not able to read the data
- 16 out of 30 documents had “broken”/“cut” tables

[REDACTED]	
Yes	Non sm Ex-drink
University engineering	
Current abode: 4 level <u>house</u> with lift Family setup: Has 3 children (2 son 1 da Main spokesperson: Eldest son Main caregiver: Maid	
Yes	

[REDACTED]	
Yes	Both; underwent right eye
Yes	Both; Hearing aids: No
Yes	Verbal - minimally <u>commu</u>





## Problems Encountered- Ren Ern

- Program extracts data by hitting certain keywords already pre-set in the codes
- Some files had empty list
- I went to open up the files to take a look, and I found out that they had different names
- Hence program was unable to catch that part of the data
- Added in that particular naming (ie, SENSORY and COMMUNI COMMUNICATION) for program to read and extract data

[illegible]

SENSORY and COMMUNICATION	
Visual Impairment	no
Hearing Impairment	yes
Communication	Yes

```
for para1 in cell.paragraphs:

    if (para1.text == 'SENSORY and COMMUNI COMMUNICATION' or para1.text == 'SENSORY and COMMUNICATION'):

        sector = 'SENSORY and COMMUNICATION'

    else:

        pass
```



# Problems Encountered- Ren Ern

- Tried to perform keyword extraction, but accuracy was too low
- Keywords were irrelevant too

```
accuracy: 27.3%
```

```
3 0.26845367061963843 grocery shopping
3 0.26845367061963843 grocery
3 0.2668648092383542 study
3 0.2668648092383542 bible study
3 0.2668648092383542 bible
3 0.23808564572147742 race indian
3 0.23808564572147742 indian sex
3 0.23479998203810487 lead
3 0.22663495651997584 sex male
3 0.21970386085598265 dob 1931
3 0.21970386085598265 1931
```



# Problems Encountered- Ren Ern

MENTAL HEALTH		
Progressive Forgetfulness or Known Cognitive impairment	No	noted Hx mild cognitive <u>cognitive</u> impairment s/b Dr Selva in 2014 noted to have STML (+) with occasional repetition and apraxia (+), (+) with occasional repetition and apraxia (+), but still able to repetition and apraxia (+), but still able to pack <u>angpao</u> apraxia (+), but still able to pack <u>angpao</u> able to pack <u>angpao</u> AMT 6/10 (no schooling) discharged from GRM <u>GRM</u>
Depression	No	Suicide risk: No looks f looks forward to weekends as grandchildren visit weekends as grandchildren visit grandchildren visit
Sleep Issues	Yes	wakes up at 8am every <u>mor</u> every morning reported she can't sleep well, wakes up x2 nightly. no <u>nocturia</u> sleep well, wakes up x2 nightly. no <u>nocturia</u> (wears diapers at x2 nightly. no <u>nocturia</u> (wears diapers at night due to falls risk) <u>nocturia</u> (wears diapers at night due to falls risk) diapers at night due to falls risk) to falls risk)
Behavioral Disturbance	No	personality: very <u>particula</u> particular, likes things to be neat and tidy things to be neat and tidy and tidy

File with comments

MENTAL HEALTH		
Progressive Forgetfulness or Known Cognitive impairment	NA	
Depression	No	
Sleep Issues	No	
Behavioral Disturbance	No	

File with no comments

IFA files are not suitable for text analysis as not every file has texts made by the doctor



## Problems Encountered- Ren Ern

SENSORY and COMMUNI COMMUNICATION	
Visual Impairment	no → 0
Hearing Impairment	yes → 1
Communication	Yes → 1

```
cat_3, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1,
cat_4, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0,
cat_4-5, 0, 0, 1, 1, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0,
cat_6-7, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
cat_2, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0,
cat_3, 0, 1, 1, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0,
cat_2, 1, 1, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0,
```

- IFA.binomial program reads Yes/No and convert it to 1/0 respectively. If it's neither Yes or No, for example, null values/missing values the program sets it to -1
- Some files were very obvious outliers, the whole dataset was set to -1



## Problems Encountered- Ren Ern

MENTAL HEALTH	
Progressive Forgetfulness or Known Cognitive impairment	Y
Depression	N
Sleep Issues	Y
Behavioral Disturbance	N



```
if (theline[0:2]=="no" or theline[0:1] == "n"):
    output += "0"
elif (theline[0:3]=="yes" or theline[0:1] == "y"):
    output += "1"
else:
    output += "-1"
```

- I went to investigate and opened up the files to take a look and some files indicated Yes/No by only stating either Y/N
- Added in some codes for program to accommodate and still convert to binomial format



## Problems Encountered- Ren Ern

SENSORY and COMMUNICATION	
Visual Impairment	No / Yes
Hearing Impairment	No / Yes
Communication	No / <b>Yes</b>



SENSORY and COMMUNICATION	
Visual Impairment	<b>No</b>
Hearing Impairment	<b>No</b>
Communication	<b>Yes</b>

- I went to investigate and opened up the files to take a look and some files indicated Yes/No by bolding the character, hence program set it as -1
- I changed the dataset myself, and only kept the bolded character



## Problems Encountered- Ren Ern

```
C:\Users\renern\Anaconda3\python.exe C:/Users/renern/Desktop/FYP/IFA_classifier.py  
the accuracy of the Naive bayes classifier on the test set is : 40.0 %  
the accuracy of the Random forest classifier on the test set is : 26.666666666666668 %  
  
Process finished with exit code 0
```

- Performance of the classifier depends heavily on the data quality
- Accuracy was too low as dataset was too little





## Problems Encountered- Ren Ern

```
the accuracy of the Naive bayes classifier on the test set is : 58.97435897435898 %  
the accuracy of the Random forest classifier on the test set is : 64.1025641025641 %  
  
Process finished with exit code 0
```

- Did oversampling to try to improve accuracy
- Accuracy was able to increase



# Problems Encountered- Ren Ern

- Predictor predicted out Category but it wasn't accurate since the accuracy of the classifier wasn't that great to begin with
- Performance of Classifier depends heavily on data quality
- We had too little data and dataset was repetitive due to oversampling
- It would not result in a fair and accurate result

The real category is : 4

The predicted category with RandomForest is : [5]

The predicted category with NaiveBayes is : [7]

The real category is : 6

The predicted category with RandomForest is : [5]

The predicted category with NaiveBayes is : [4]

The real category is : 7

The predicted category with RandomForest is : [5]

The predicted category with NaiveBayes is : [4]

# Summing-up - Ren Ern

**IFA\_json\_extract.py**  
(for extracting data out from IFA files)

A method was created in the extract to run it in the Classifier when called



**IFA\_binomial.py**  
(for converting Yes/Nos to 1/0s)

A method was created in the script to run it in the Classifier when called



**IFA\_classifier.py**  
(to build a model and predict categories)

Classifier runs for both **IFA\_json\_extract.py** and **IFA\_binomial.py** when program is run

Q & A





**Thank You for Listening!**