



Mahidol University

Faculty of Medicine Ramathibodi Hospital

Department of Clinical Epidemiology and Biostatistics

Using Natural Language Processing (NLP) to predict colon cancer from discharge summaries

CEB workshop

25 November 2022 13:00 – 15:00

Wanchana Ponthongmak

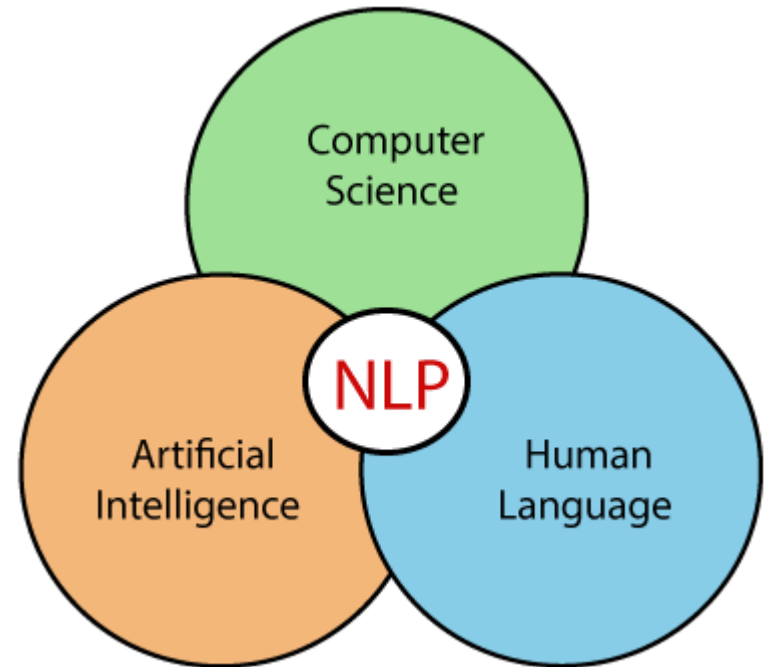
Presentation available here!

shorturl.at/eV348



What is Natural Language Processing (NLP)?

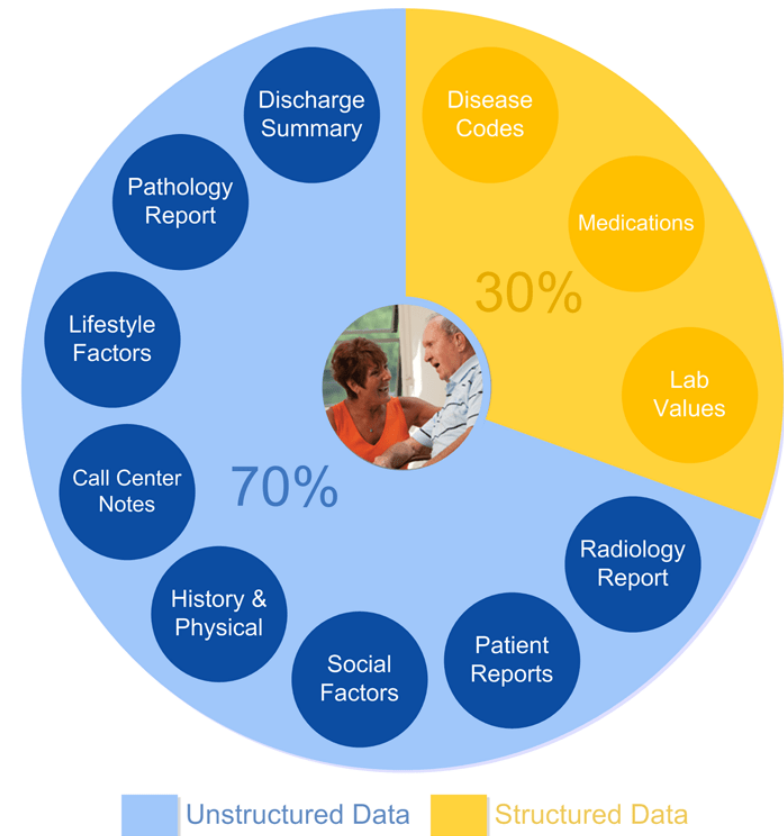
- NLP is a multi-disciplinary fields, which is a part of computer science, human language, and artificial intelligence (AI).
- It offers computers ability to analyze, understand, and utilize human languages as they are spoken and written,
 - Textual form
 - Vocal form





Why NLP?

- Text data are extremely rich source of information
- But extracting insights from them can be hard and time-consuming, due to its unstructured nature.
- ~70% of data in hospitals are unstructured data





Mahidol University

Faculty of Medicine Ramathibodi Hospital

Department of Clinical Epidemiology and Biostatistics

Why NLP?

- But, thanks to advances in NLP, machine learning (ML), deep learning (DL) and AI, sorting text data is getting easier.



<https://online.york.ac.uk/the-role-of-natural-language-processing-in-ai/>



Mahidol University

Faculty of Medicine Ramathibodi Hospital

Department of Clinical Epidemiology and Biostatistics

Why NLP?

NLP Applications

Frequent applications of NLP are



<https://www.forbesindia.com/article/weschool/natural-language-processing-a-breakthrough-technology-in-healthcare/67661/1>



Mahidol University

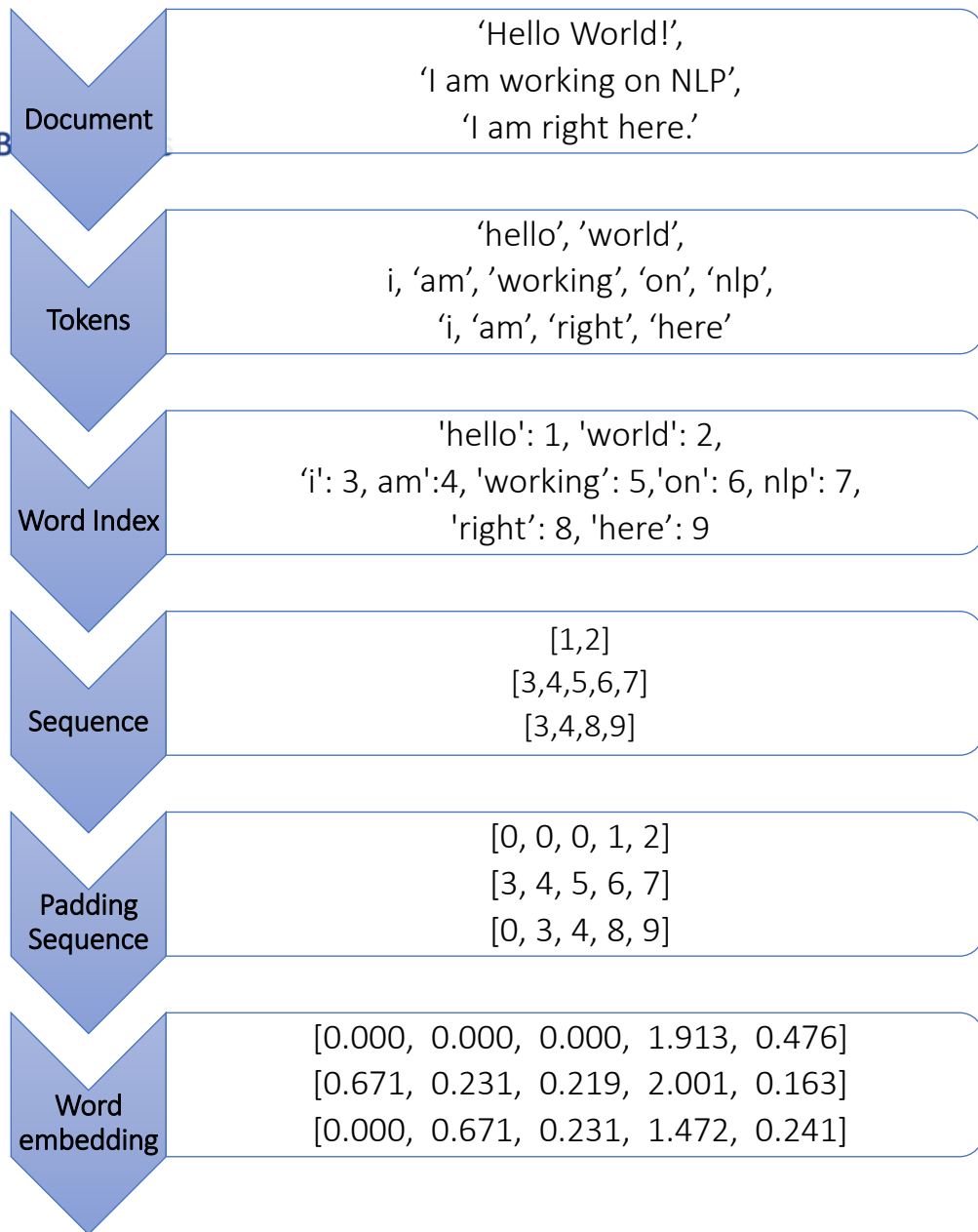
Faculty of Medicine Ramathibodi Hospital

Department of Clinical Epidemiology and Biostatistics

How to deal with text data?



- Structurize text data
 - vectorize to numerical vector with the same size





Mahidol University

Faculty of Medicine Ramathibodi Hospital

Department of Clinical Epidemiology and Biostatistics

Objectives of the workshop

1

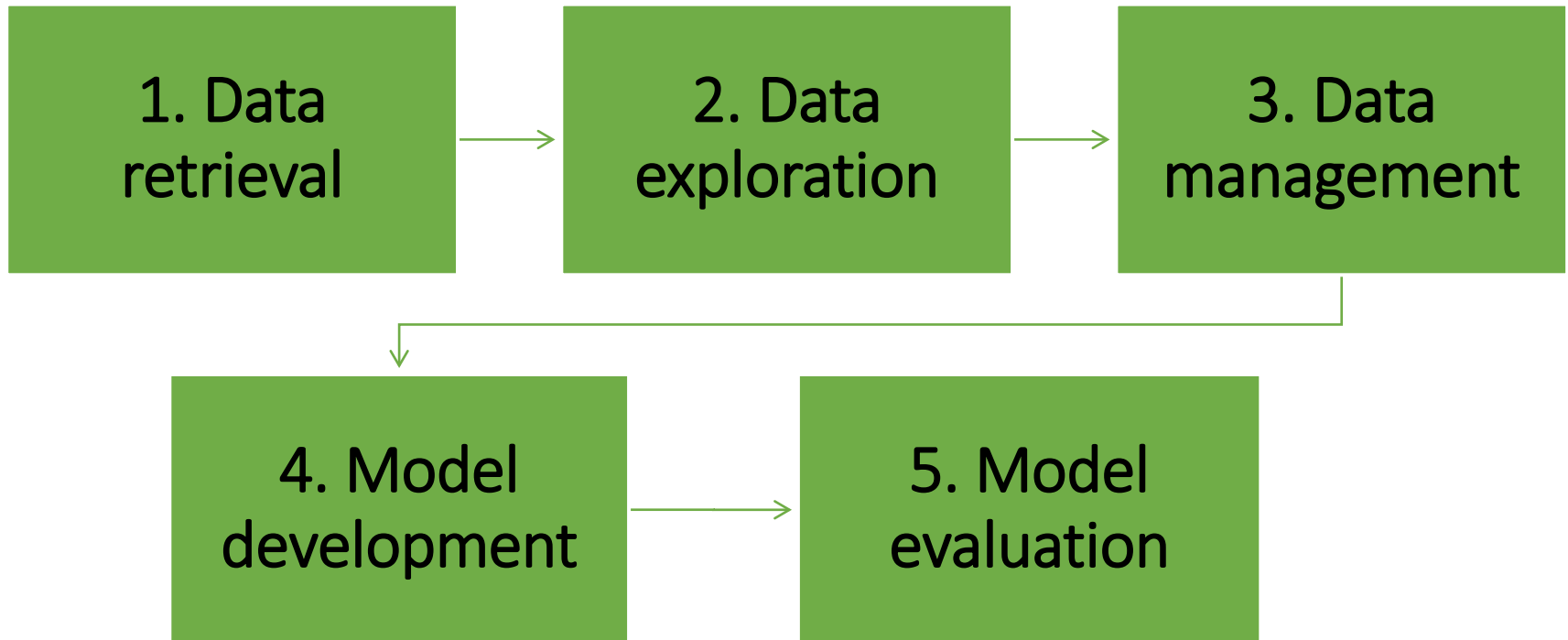
To develop models to predict colon cancer from discharge summary using NLP with machine learning

2

To evaluate the performance of developed models



Framework





Mahidol University

Faculty of Medicine Ramathibodi Hospital

Department of Clinical Epidemiology and Biostatistics

Data Retrieval

Study design and setting



A cross-sectional study from

Department of Medicine, Ramathibodi Hospital

1st January 2015 – 31st December 2019,

Ramathibodi Human Research Ethics Committee
(COA.MURA2020/1152).

Inclusions

- Audited charts
- ≥ 18 years old
- Top 10th diseases

Exclusions

- Missing diagnosis codes
- Missing clinical notes

1. Data retrieval

2. Data exploration

3. Data management

4. Model
development

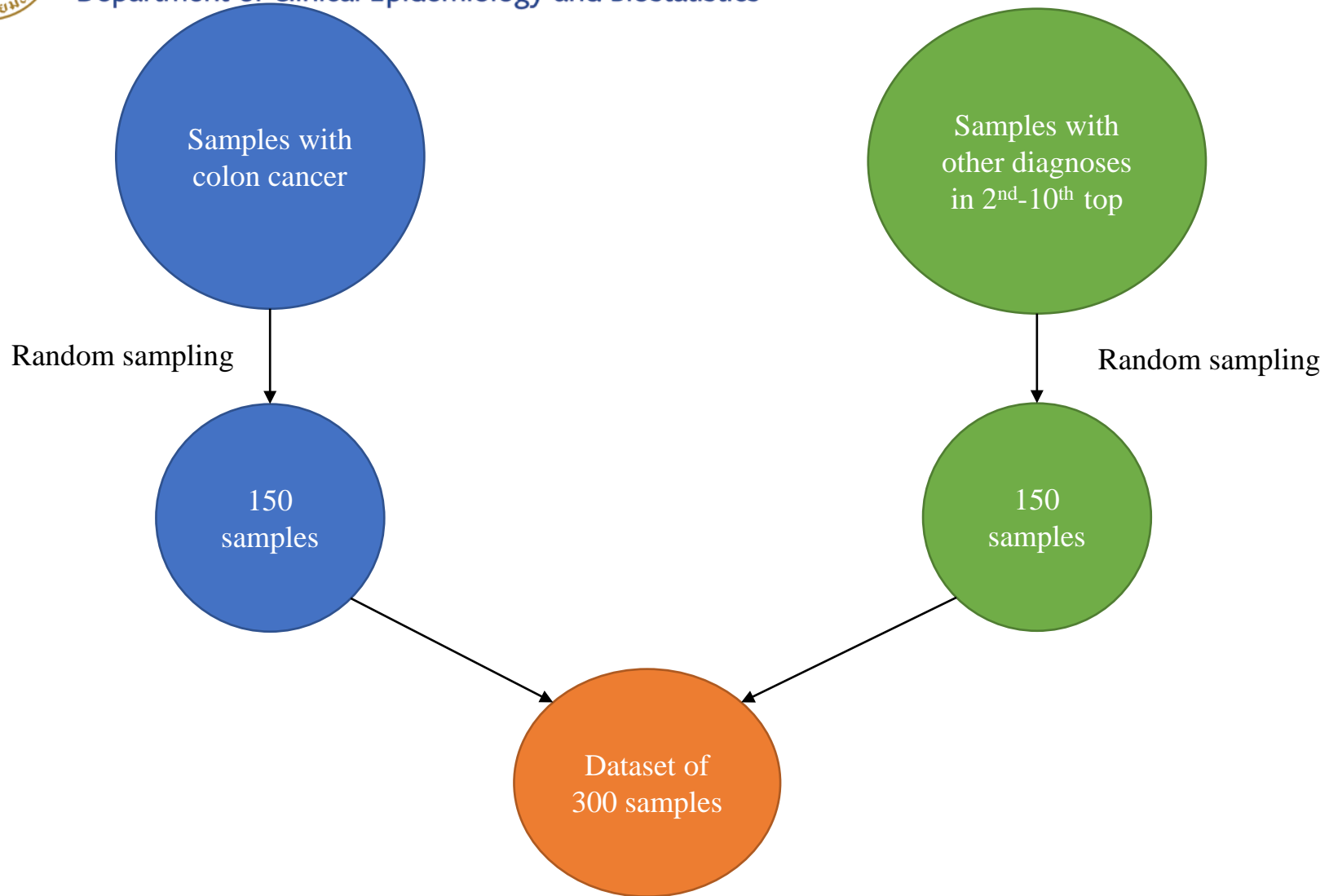
5. Model evaluation



Mahidol University

Faculty of Medicine Ramathibodi Hospital

Department of Clinical Epidemiology and Biostatistics



1. Data retrieval

2. Data exploration

3. Data management

4. Model
development

5. Model evaluation



Data Exploration

The data exploration applied descriptive statistics and data visualization techniques to understand the characteristic of the study datasets

- Features exploration

- Data characteristics

- # of patients
- # of documents
- Age / gender
- # of colon CA
- Vocabulary size
- # of tokens per documents

hn	an	age	gender	pdx	brief_raw	course_raw	label
wpLCImZowqrCo	wpLCmGpmwq3C	49	2	C20	known case CA rectum Dx. l	CA rectosigmoid colon c ova	1
wpPCmmNqwqv	wpLCmWRswq3C	67	1	C20	case CA lower rectum S/P A	#recurrent CA lower rectum	1
wpXCk2lpwrHCpl	wpLCmWtswrHC	45	2	C20	known case ca rectum c live	# CA rectal with liver metas	1
wpbCl2pswq7Cpl	wpLCm2RswqzCr	81	2	C20	Known case CA Rectum with	case CA Rectum with liver n	1
wpbCl2puwqzCn	wpLCm2Vuwq_C	62	2	C20	Known case : CA rectum (T3	S/P mFOLFOX6 1st cycle 3-4,	1
wpbCk2huwrLC	wpLCmWhswqrC	67	2	C20	Known case : CA upper rect	1. Recurrent CA rectum 151	1
wpbCl2NnwqvCr	wpLCmmxlwrLC	71	1	C20	known case CA rectum S/P L	#CA rectum S/P LAR start m	1
wpPCI2xuwrlCo	wpLCmmlswrDC	57	1	C20	#Known case CA rectum pre	#Known case CA rectum T3N	1
wpbCk2duwq7Cr	wpLCmWhmwrLC	65	2	C20	Known case CA lower rectu	# CA lower rectum with LN r	1
wplCkmdswqnC	wpLCmGlnwrHC	68	2	C20	Known case # CA rectum wi	# CA rectum with liver meta	1
wpxCm2luwq_Cr	wpLCmGtnwrHC	56	1	C20	#U/D DM, HT, old CVA + yrs	#U/D DM, HT, old CVA 10+ y	1
wpbClmlqwq7Co	wpLCmmtrwqyCo	70	1	C20	U/D IHD S/P CAG 23/11/61	#CA rectosigmoid pT3N1M0	1
wpxCk2ptwqnCo	wpLCmmtlwq7Cr	76	1	C20	Known case # DLP # small bc	# DLP # small bowel GIST s/	1
wpbCl2Rmwq7C	wpLCmmtrwqzCo	57	1	C20	Known case CA upper rectu	#CA upper rectum with mul	1
wpbClWpowq7C	wpLCmmpuwrDC	46	1	C20	known case CA lower rectu	#CA lower rectum stage IIIB	1

Textual data

Sign to save

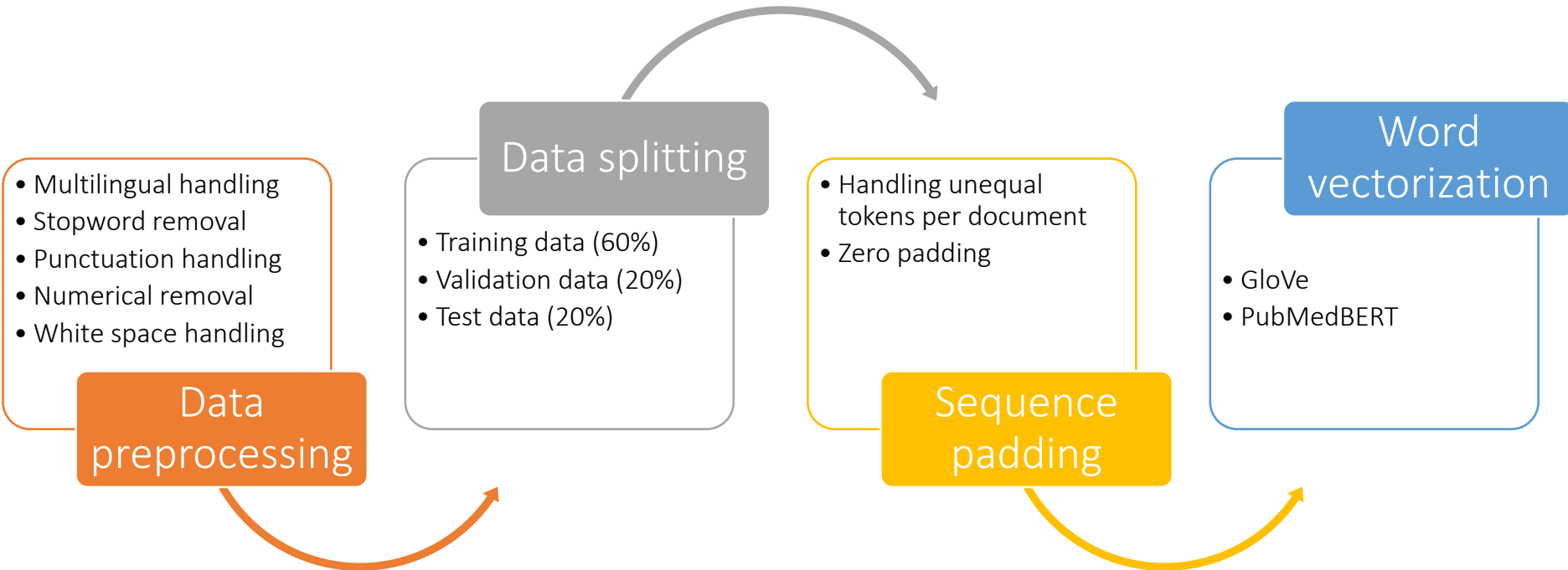
12 NW 2563

RH.0011

5. Model evaluation



Data Management





Mahidol University

Faculty of Medicine Ramathibodi Hospital

Department of Clinical Epidemiology and Biostatistics

Model Development

Two models were applied

1. Gated Recurrent Unit (GRU)
 - Pre-trained GloVe (Global Vectors for Word Representation)
2. Bidirectional Encoder Representations from Transformers (BERT)
 - Pre-trained BERT on a medical domain (PubMedBERT)

1. Data retrieval

2. Data exploration

3. Data management

4. Model
development

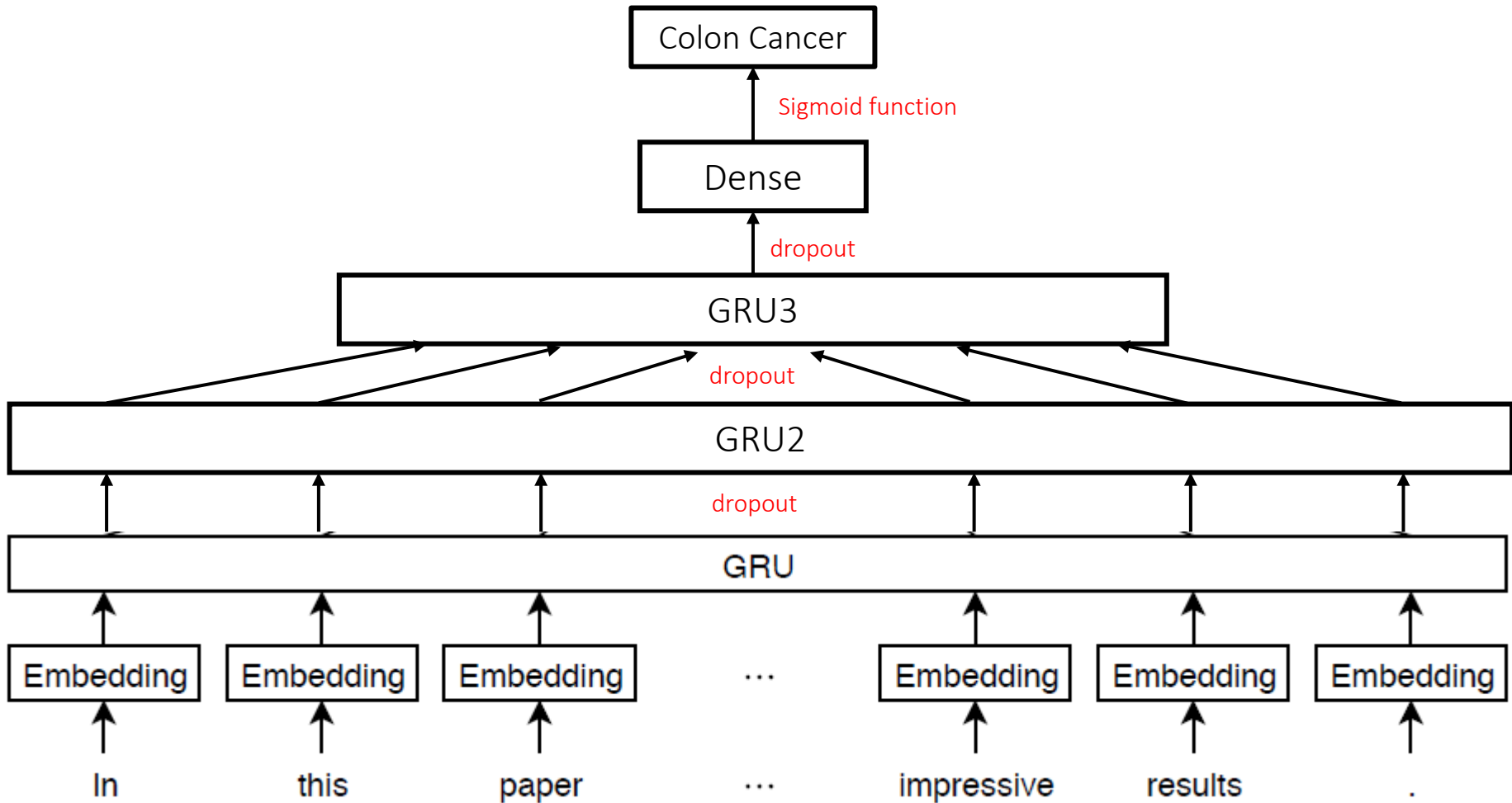
5. Model evaluation



Mahidol University

Faculty of Medicine Ramathibodi Hospital

Department of Clinical Epidemiology and Biostatistics



1. Data retrieval

2. Data exploration

3. Data management

4. Model
development

5. Model evaluation



Mahidol University

Faculty of Medicine Ramathibodi Hospital

Department of Clinical Epidemiology and Biostatistics

Layer (type)	Output Shape	Param #
embedding (Embedding)	(None, 241, 50)	253700
batch_normalization (Batch Normalization)	(None, 241, 50)	200
dropout (Dropout)	(None, 241, 50)	0
gru (GRU)	(None, 241, 32)	8064
dropout_1 (Dropout)	(None, 241, 32)	0
gru_1 (GRU)	(None, 241, 32)	6336
dropout_2 (Dropout)	(None, 241, 32)	0
gru_2 (GRU)	(None, 16)	2400
dropout_3 (Dropout)	(None, 16)	0
dense (Dense)	(None, 1)	17
Total params: 270,717		
Trainable params: 16,917		
Non-trainable params: 253,800		

1. Data retrieval

2. Data exploration

3. Data management

4. Model
development

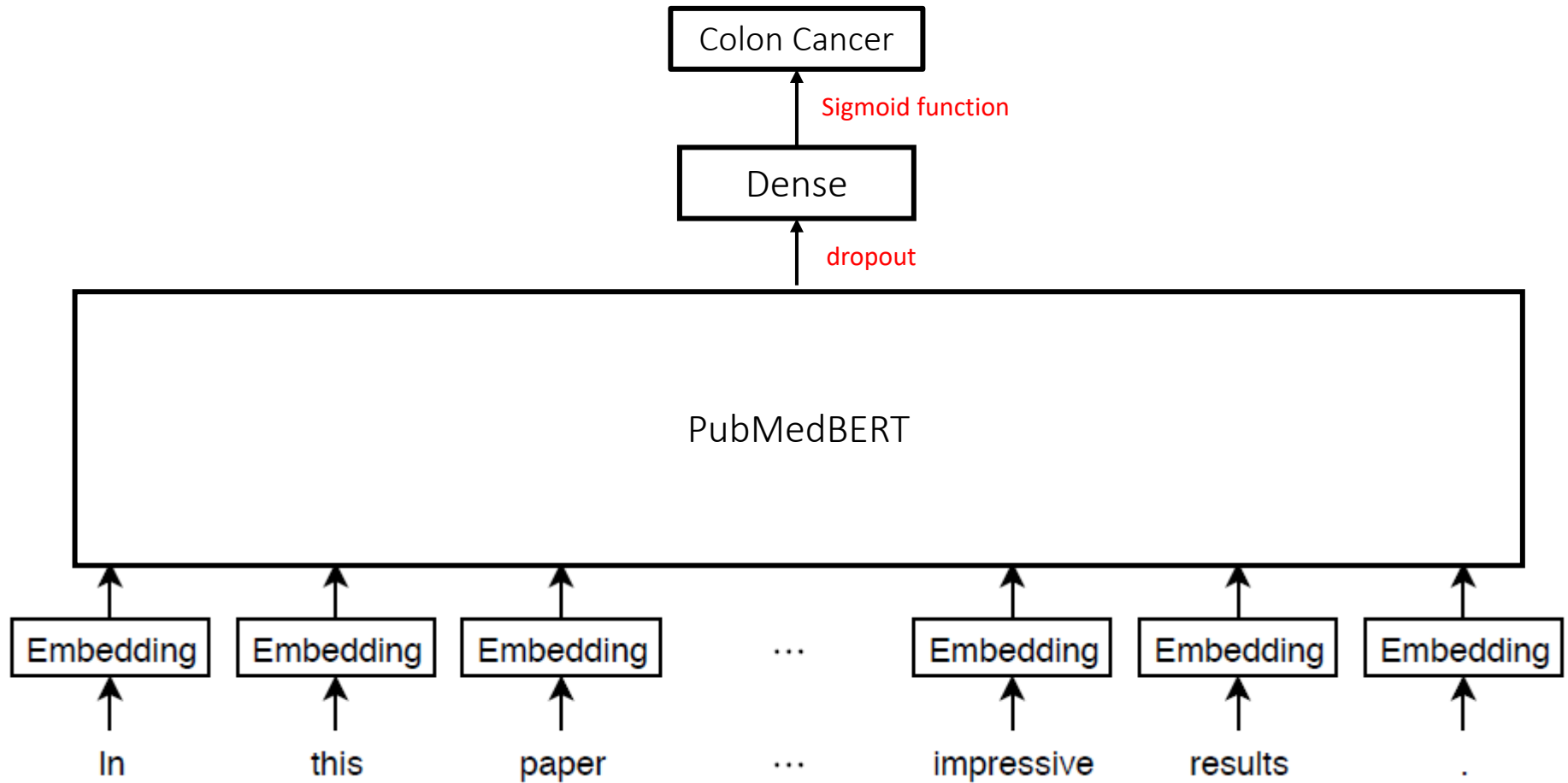
5. Model evaluation



Mahidol University

Faculty of Medicine Ramathibodi Hospital

Department of Clinical Epidemiology and Biostatistics



1. Data retrieval

2. Data exploration

3. Data management

4. Model
development

5. Model evaluation



M
Fa
De

Layer (type)	Output Shape	Param #	Connected to
text (InputLayer)	[(None,)]	0	[]
keras_layer (KerasLayer)	{'input_mask': (None, 128), 'input_word_ids': (None, 128), 'input_type_ids': (None, 128)}	0	['text[0][0]']
keras_layer_1 (KerasLayer)	{'sequence_output': (None, 128, 768), 'encoder_outputs': [(None, 128, 768), (None, 128, 768), (None, 128, 768), (None, 128, 768), (None, 128, 768), (None, 128, 768), (None, 128, 768), (None, 128, 768), (None, 128, 768), (None, 128, 768), (None, 128, 768), (None, 128, 768)], 'pooled_output': (None, 768), 'default': (None, 768)}	109482241	['keras_layer[0][0]', 'keras_layer[0][1]', 'keras_layer[0][2]']
dropout (Dropout)	(None, 768)	0	['keras_layer_1[0][13]']
output (Dense)	(None, 1)	769	['dropout[0][0]']
Total params: 109,483,010 Trainable params: 769 Non-trainable params: 109,482,241			

1. Data retrieval

2. Data exploration

3. Data management

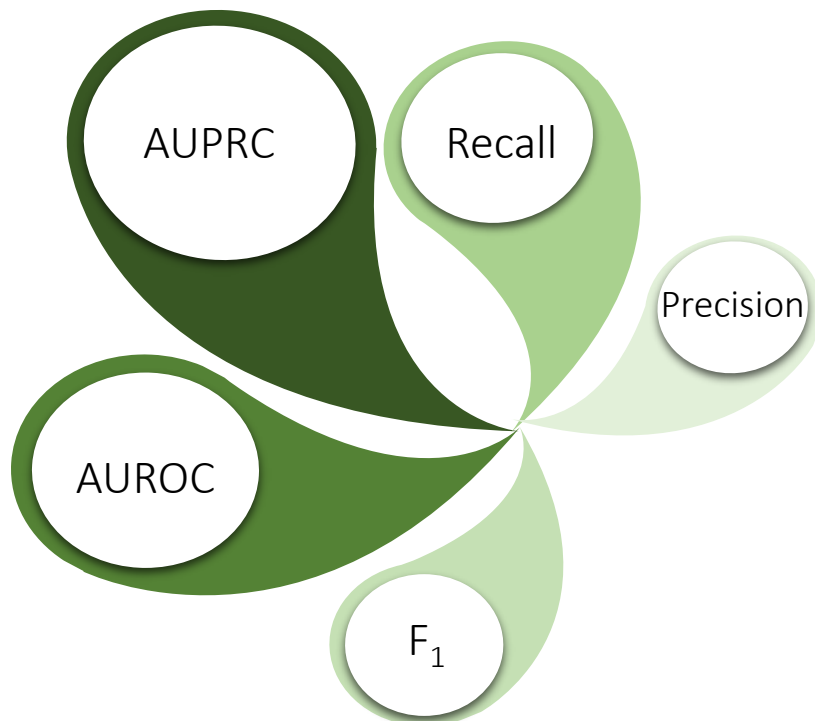
4. Model
development

5. Model evaluation



Model evaluation

Evaluation metrics



Monitor model overfitting by the percent difference between training and test performances

1. Data retrieval

2. Data exploration

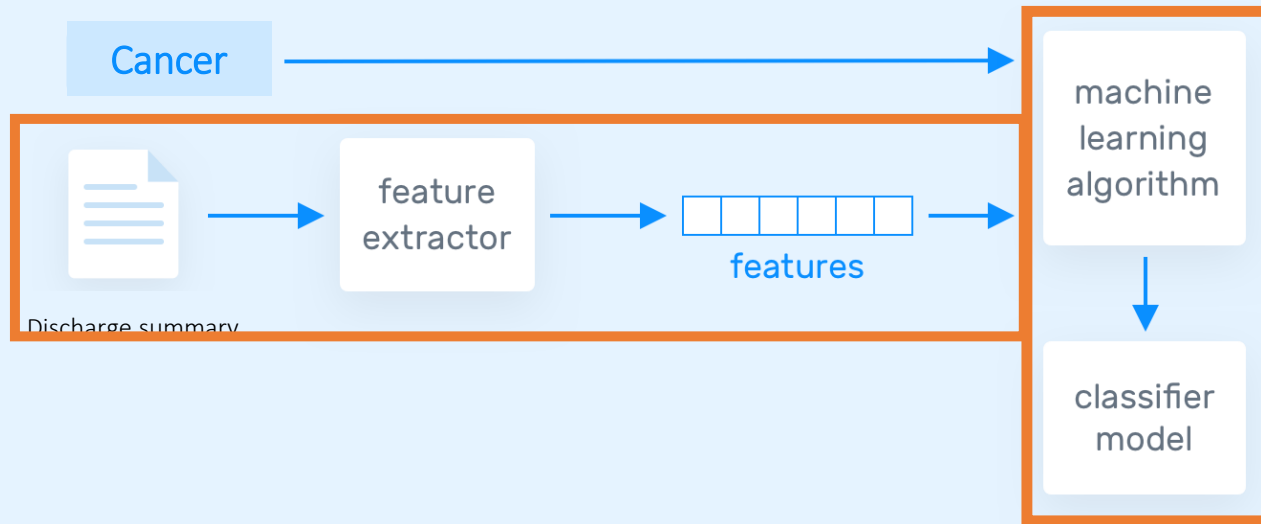
3. Data management

4. Model development

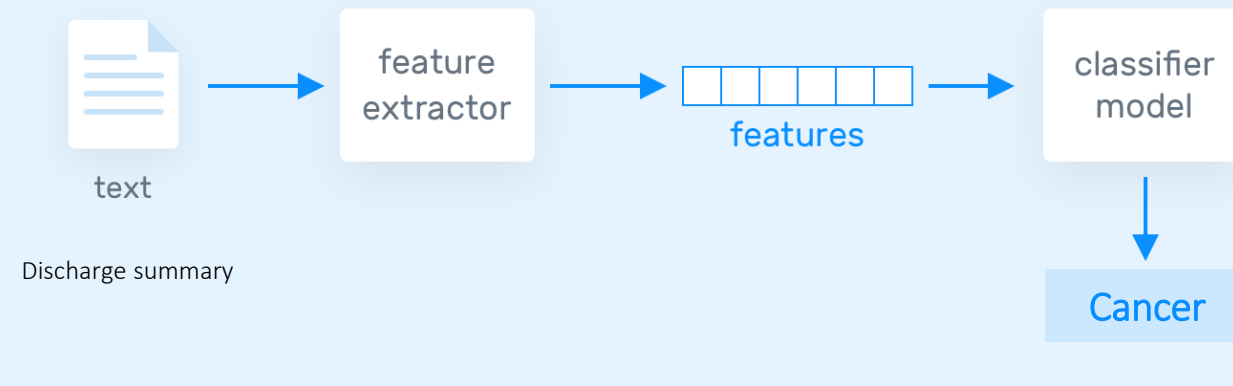
5. Model evaluation



(a) Training



(b) Prediction



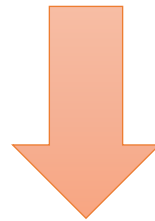


Mahidol University

Faculty of Medicine Ramathibodi Hospital

Department of Clinical Epidemiology and Biostatistics

Let's get
your hands dirty



shorturl.at/mql24