

การสรุปข้อความข่าวภาษาไทย
โดยใช้แบบจำลองการเรียนรู้เชิงลึก

Thai News Text Summarization
Using Deep Learning Model



ปัญหาพิเศษนี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตร
ปริญญาวิทยาศาสตรบัณฑิต (วิทยาการคอมพิวเตอร์)
ภาควิชาวิทยาการคอมพิวเตอร์ คณะวิทยาศาสตร์
สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง
ปีการศึกษา 2561

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Thai News Text Summarization

Using Deep Learning Model



A SPECIAL PROBLEM SUBMITTED IN PARTIAL FULFILLMENT OF
THE REQUIREMENTS FOR
THE DEGREE OF BACHELOR OF SCIENCE (COMPUTER SCIENCE)
DEPARTMENT OF COMPUTER SCIENCE, FACULTY OF SCIENCE
KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG
ACADEMIC YEAR 2018

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

หัวข้อปัญหาพิเศษ การสรุปข้อความข่าวภาษาไทยโดยใช้แบบจำลองการเรียนรู้เชิงลึก
Thai News Text Summarization Using Deep Learning Model
ชื่อนักศึกษา นายธีรพงศ์ แต่งเมือง รหัสนักศึกษา 58050290
ปริญญา วิทยาศาสตร์บัณฑิต
ภาควิชา วิทยาการคอมพิวเตอร์
อาจารย์ที่ปรึกษา ดร.กุลสวัสดิ์ จิตขจรวานิช

คณะวิทยาศาสตร์ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง (สจล.)
อนุมัติให้ปัญหาพิเศษนี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิทยาศาสตรบัณฑิต
(วิทยาการคอมพิวเตอร์) ประจำปีการศึกษา 2561

คณะกรรมการสอบ	ลายมือชื่อ
ผศ.ดร.สายชล ใจเย็น ประธานกรรมการ	
ผศ.ดร.อนันตพร หารรรษคุณาตย์ กรรมการ	
ดร.กุลสวัสดิ์ จิตขจรวานิช กรรมการและอาจารย์ที่ปรึกษา	

ลิขสิทธิ์ของคณะวิทยาศาสตร์
สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

หัวข้อปัญหาพิเศษ	การสรุปข้อความข่าวภาษาไทยโดยใช้แบบจำลองการเรียนรู้เชิงลึก
ชื่อนักศึกษา	นายธีรพงศ์ แต่งเมือง รหัสนักศึกษา 58050290
ปริญญา	วิทยาศาสตร์บัณฑิต
ภาควิชา	วิทยาการคอมพิวเตอร์
คณะ	วิทยาศาสตร์
มหาวิทยาลัย	สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง (สจล.)
ปีการศึกษา	2561
อาจารย์ที่ปรึกษา	ดร.กุลสวัสดิ์ จิตขจรวานิช

บทคัดย่อ

เนื่องด้วยการศึกษาและพัฒนาเทคนิคการสรุปข้อความสำหรับภาษาไทยยังมีข้อจำกัด ปัญหาพิเศษนี้มีวัตถุประสงค์เพื่อแก้ไขปัญหาดังกล่าวโดยการประยุกต์ใช้เทคนิคการเรียนรู้เชิงลึกซึ่งเป็นศาสตร์ที่กำลังได้รับความนิยม ข้อมูลที่นำมาใช้ในการศึกษาค้นคว้าครั้งนี้ได้รวบรวมมาจากข้อมูลข่าวจากสถานีโทรทัศน์ไทยพีบีเอสโดยเลือกเฉพาะส่วนของหัวข้อข่าวและข่าวย่อ จำนวน 29,170 ข่าว มาใช้ในการสร้างแบบจำลองการเรียนรู้ของเครื่อง จากการดำเนินงานพบว่าเป็นการเริ่มต้นศึกษาและมีแนวโน้มในการพัฒนาไปในทางที่ดี อย่างไรก็ตามผลการวัดประสิทธิภาพด้วยตัวชี้วัด ROUGE พบว่ายังมีความถูกต้องน้อย ซึ่งคาดว่าต้องใช้ข้อมูลข่าวอีกจำนวนมากในการสร้างแบบจำลองการเรียนรู้เชิงลึก อีกทั้งยังต้องมีการศึกษาและทดลองนำอัลกอริทึมใหม่ๆ มาประยุกต์ใช้กับการสรุปข้อความภาษาไทยอีกจำนวนมากเพื่อให้ได้ผลลัพธ์ที่ดียิ่งขึ้น

คำสำคัญ: การประมวลผลภาษาธรรมชาติ การเรียนรู้ของเครื่อง การเรียนรู้เชิงลึก การสรุปข้อความ การสรุปข่าว

Title	Thai News Text Summarization Using Deep Learning Model
Students	Mr. Teerapong Taengmuang Student ID 58050290
Degree	Bachelor of Science
Department	Computer Science
Faculty	Science
University	King Mongkut's Institute of Technology Ladkrabang (KMITL)
Academic Year	2018
Advisor	Dr. Kulsawasd Jitkajornwanich

Abstract

Designing and developing an efficient text summarization tool for Thai language is still considered to be a challenge. The goal of this special problem is to overcome such challenge by applying deep learning technique, one of the increasingly popular fields of study in the era of big data. The Thai news dataset used in this work is collected from Thai PBS TV station, consisting of 29,170 news in which only news headlines and short news were used in developing a machine learning model. The results showed that the proposed model has a good potential evidenced by a performance evaluation, called ROUGE. However, there is still room for improvements especially for the accuracy aspect. In future work, we expect that with more and larger news datasets, our deep learning-based model will perform better. In addition, applying more diverse/new algorithms of text summarization should also be considered in improving the overall performance of the Thai news text summarization system.

Keywords: Natural Language Processing, Machine Learning, Deep Learning, Text Summarization, News Summarization

กิตติกรรมประกาศ

โครงการพิเศษฉบับนี้สำเร็จลุล่วงไปได้ด้วยดี เนื่องจากความช่วยเหลือและการสนับสนุนจากบุคคลหลายๆ ท่าน ขอขอบคุณอาจารย์ที่ปรึกษา ดร.กุลสวัสดิ์ จิตขจรวานิช ผู้ซึ่งให้คำปรึกษาและคำแนะนำอันเป็นประโยชน์ต่อการปรับปรุงข้อบกพร่องในการจัดทำโครงการพิเศษ ขอขอบคุณ ผศ.ดร.สายชล ใจเย็น และ ผศ.ดร.อนันตพร ทรรษคุณาตม์ ประธานกรรมการและกรรมการ ที่ให้ข้อคิดเห็นและข้อเสนอแนะในการจัดทำโครงการพิเศษให้สำเร็จลุล่วงไปได้ด้วยดี

ขอกราบขอบพระคุณครูบาอาจารย์ทุกท่านที่ประสิทธิ์ประสาทวิชาความรู้ให้แก่ข้าพเจ้าตั้งแต่เยาว์วัยจนข้าพเจ้ามีความรู้ความสามารถในการจัดทำโครงการพิเศษนี้ได้

และสุดท้ายข้าพเจ้าขอกราบขอบพระคุณ บิดา มารดา และบุคคลในครอบครัว ที่ให้โอกาสในการศึกษาอย่างเต็มที่ ตลอดจนเลี้ยงดู อบรมสั่งสอน พร้อมทั้งให้กำลังใจและความรักเสมอมา

ธีรพงศ์ แต่งเมือง

สารบัญ

	หน้า
บทคัดย่อภาษาไทย	ก
บทคัดย่อภาษาไทย	ข
กิตติกรรมประกาศ	ค
สารบัญ	ง
สารบัญตาราง	ช
สารบัญรูป	ณ
บทที่ 1 บทนำ	1
1.1 ที่มาและความสำคัญของปัญหาพิเศษ	1
1.2 วัตถุประสงค์ของปัญหาพิเศษ	1
1.3 ขอบเขตของปัญหาพิเศษ	2
1.4 ประโยชน์ที่คาดว่าจะได้รับ	2
1.5 ขั้นตอนการดำเนินงาน	2
1.6 อุปกรณ์ที่ใช้ในการทำปัญหาพิเศษ	3
บทที่ 2 ทฤษฎีและเอกสารที่เกี่ยวข้อง	4
2.1 การตัดคำ	4
2.1.1 วิธีการตัดคำ	5
2.1.2 เทคนิคการตัดคำ	6
2.2 Word2Vec	7
2.3 ขนาดของแบตช์	9
2.4 Epoch	9
2.5 อัตราการเรียนรู้	9
2.6 การลดอัตราการเรียนรู้	9
2.7 One-Hot Encoding	10
2.8 ฟังก์ชันความสูญเสีย	10
2.9 Categorical Cross-Entropy	11

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญ (ต่อ)

	หน้า
2.10 Optimizer	12
2.11 ROUGE	12
2.11.1 ROUGE-1	13
2.11.2 ROUGE-2	13
2.11.3 ROUGE-L	14
2.12 F-measure	14
2.13 Longest Common Subsequence Problem	15
2.14 การเรียนรู้เชิงลึก	19
2.15 การเข้ากันมากเกินไป	19
2.16 Dropout	20
2.17 การเรียนแบบวนซ้ำ	21
2.17.1 ปัญหาของ RNN	23
2.18 Long Short-Term Memory (LSTM)	23
2.18.1 การทำงานของ LSTM	24
2.19 Bi-directional LSTM	26
2.19.1 ขั้นตอนการทำงาน Bi-directional LSTM	27
2.20 แบบจำลอง Sequence-to-Sequence	29
2.21 Bahdanau Attention	30
2.22 เอกสารและงานวิจัยที่เกี่ยวข้อง	35
2.22.1 เอกสารที่เกี่ยวข้อง	35
2.22.2 งานวิจัยที่เกี่ยวข้อง	37

สารบัญ (ต่อ)

	หน้า
บทที่ 3 วิธีการดำเนินงาน	38
3.1 กระบวนการจัดเตรียมข้อมูล	38
3.1.1 การเก็บรวบรวมข้อมูล	39
3.1.2 แปลงคำศัพท์ย่อ	40
3.1.3 ตัดคำ	40
3.1.4 กำจัดคำหยุด	43
3.1.5 กำจัดเครื่องหมายสัญลักษณ์	43
3.1.6 แปลงคำที่ไม่มีในพจนานุกรมเวกเตอร์เป็น “<UNK>”	43
3.1.7 เพิ่ม “<PAD>”, “<EOS>”, “<GO>”	44
3.1.8 แปลงคำศัพท์เป็นดัชนีของคำศัพท์	45
3.1.9 แปลงดัชนีของคำศัพท์เป็นค่าเวกเตอร์ของคำศัพท์	46
3.2 สร้างแบบจำลองการเรียนรู้ของเครื่อง	46
3.2.1 การกำหนดค่า Hyperparameter	46
3.2.2 การออกแบบจำลองการเรียนรู้ของเครื่อง	47
3.3 การวัดประสิทธิภาพของแบบจำลองการเรียนรู้ของเครื่อง	48
บทที่ 4 ผลการทดลองและอภิปรายผล	49
4.1 ผลการทดลอง	49
4.2 อภิปรายผลการทดลอง	52
บทที่ 5 สรุปผลการทดลองและข้อเสนอแนะ	53
5.1 สรุปผลการทดลอง	53
5.2 ข้อเสนอแนะ	53

สารบัญ (ต่อ)

	หน้า
เอกสารอ้างอิง	54
ภาคผนวก	57
ภาคผนวก ก รายการคำศัพท์ย่อ	58
ภาคผนวก ข ตัวอย่างผลลัพธ์การสรุปข้อความข่าว	59
ภาคผนวก ค โค้ดที่ใช้ในการสร้างแบบจำลอง	62



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญตาราง

ตารางที่	หน้า
2.1 เกณฑ์การตัดคำภาษาไทย	4
2.2 การเปรียบเทียบข้อดีและข้อเสียของวิธีการตัดคำภาษาไทย	5
2.3 เริ่มต้นสร้างตาราง LCS	15
2.4 การคำนวณ LCS กับข้อมูลแถวที่ 2	16
2.5 การคำนวณ LCS กับข้อมูลแถวที่ 3	17
2.6 การคำนวณ LCS กับข้อมูลแถวที่ 4	17
2.7 จำนวนความยาวสูงสุดของชุดอักษรในแถวของตาราง	18
2.8 แสดงผลลัพธ์การทำ Dot Product ระหว่าง Decoder Hidden State และ Encoder Hidden State	32
2.9 ผลลัพธ์จากการผ่านฟังก์ชัน Softmax	33
2.10 ผลลัพธ์จากการคูณ Encoder Hidden State และ Score'	34
4.1 แสดงผลลัพธ์ค่าความถูกต้องโดยตัวชี้วัด ROUGE-1, ROUGE-2 และ ROUGE-L ประกอบด้วยค่าชี้วัด Recall, Precision และ F-measure	51

สารบัญรูป

รูปที่	หน้า
2.1 แสดงความสัมพันธ์ของคำ	8
2.2 แสดงการลดลงของอัตราการเรียนรู้	10
2.3 โครงข่ายประสาทเทียมแบบลึก	19
2.4 แสดงการ Dropout ของแบบจำลองโครงข่ายประสาทเทียม	20
2.5 สถาปัตยกรรม RNN	21
2.6 สถาปัตยกรรม LSTM	23
2.7 สถาปัตยกรรม Bi-directional LSTM	27
2.8 Bi-directional LSTM ขั้นตอนไปข้างหน้า	27
2.9 Bi-directional LSTM ขั้นตอนย้อนกลับ	28
2.10 การ Concatenate	28
2.11 แบบจำลอง Seq2Seq	29
2.12 Attention Layer	30
2.13 Attention Layer การคำนวณขั้นตอนที่ 1	31
2.14 Attention Layer การคำนวณขั้นตอนที่ 2	31
2.15 Attention Layer การคำนวณขั้นตอนที่ 3	32
2.16 Attention Layer การคำนวณขั้นตอนที่ 4	33
2.17 Attention Layer การคำนวณขั้นตอนที่ 5	34
2.18 Attention Layer การคำนวณขั้นตอนที่ 6	35
3.1 เว็บไซต์ ThaiPBS	39
3.2 ส่วนของข้อมูลข่าวที่นำมาประมวลผล	40
3.3 ขั้นตอนการดำเนินการกำหนดขอบเขตคำด้วยเทคนิค newmm กับคำว่า “ฮักเจ้าหลาย”	41
3.4 แสดงอัตราส่วนของคำที่พบและไม่พบในพจนานุกรมเวกเตอร์	43
3.5 แสดงการลดลงของอัตราการเรียนรู้ที่คำนวณได้จากสมการ 3.1 เมื่อกำหนดให้อัตราการเรียนรู้เริ่มต้นเท่ากับ 0.005	47
3.6 แบบจำลอง Bahdanau	48
4.1 แสดงค่าความสูญเสียของแบบจำลองในแต่ละ Epoch	50

บทที่ 1

บทนำ

1.1 ที่มาและความสำคัญของปัญหาพิเศษ

ในปัจจุบันมีข้อมูลในรูปแบบข้อความจำนวนมากซึ่งผู้อ่านต้องใช้เวลาในการอ่านและทำความเข้าใจกับเนื้อหา แต่ในบางโอกาสผู้อ่านเพียงต้องการทราบถึงใจความสำคัญหรือสาระสำคัญเท่านั้น

จากปัญหาดังกล่าวจึงมีการพยายามพัฒนาเครื่องมือและวิธีการสำหรับสรุปข้อความตลอดมา เพื่อลดเวลาในการอ่านและสามารถเข้าใจถึงสาระสำคัญได้อย่างรวดเร็ว

การเรียนรู้เชิงลึก (Deep Learning) เป็นสาขาวิชาที่ได้รับความนิยมมากขึ้นเรื่อย ๆ เนื่องจากสามารถแก้ปัญหาที่วิธีการอื่นไม่สามารถแก้ปัญหาคือหรือให้ความถูกต้องมากกว่าวิธีการอื่นในบางปัญหา ด้วยความสามารถดังกล่าวทำให้แบบจำลองการเรียนรู้เชิงลึกได้รับการศึกษาและพัฒนาอย่างแพร่หลาย

ในปัญหาพิเศษนี้ได้นำแบบจำลองการเรียนรู้เชิงลึกสำหรับการประมวลผลธรรมชาติ (Natural Language Processing) คือแบบจำลอง Sequence-to-Sequence (Seq2Seq) มาปรับใช้กับข้อความข่าวภาษาไทยเพื่อสรุปข้อความข่าวให้สั้นลงและเป็นแนวทางในการพัฒนาการสรุปข้อความในภาษาไทยต่อไป

1.2 วัตถุประสงค์ของปัญหาพิเศษ

- 1) เพื่อลดเวลาในการอ่านข่าว
- 2) เพื่อทดลองนำแบบจำลองประเภท Sequence-to-Sequence มาปรับใช้กับการสรุปข้อความข่าวภาษาไทย
- 3) เพื่อพัฒนาองค์ความรู้ในสาขาวิชาการประมวลผลภาษาธรรมชาติในภาษาไทย

1.3 ขอบเขตของปัญหาพิเศษ

- 1) ใช้ Bahdanau ในการสร้างแบบจำลองการเรียนรู้ของเครื่อง
- 2) ใช้ Word2Vec ในการสกัดคุณลักษณะจากข้อความข่าว
- 3) รวบรวมข่าวจากสถานีโทรทัศน์ไทยพีบีเอส (Thai Public Broadcasting Service: ThaiPBS)
- 4) ตัวตัดคำภาษาไทย คือ newmm

1.4 ประโยชน์ที่คาดว่าจะได้รับ

- 1) ลดเวลาในการอ่านข่าว
- 2) เป็นแนวทางในการศึกษาการสรุปข้อความข่าวภาษาไทยโดยใช้แบบจำลองการเรียนรู้เชิงลึกและแบบจำลองประเภท Sequence-to-Sequence
- 3) พัฒนาคณะความรู้ในสาขาวิชาการประมวลผลภาษาธรรมชาติในภาษาไทย

1.5 ขั้นตอนการดำเนินงาน

- 1) ศึกษางานวิจัยและทฤษฎีที่เกี่ยวข้อง
- 2) ศึกษาการสกัดคุณลักษณะด้วยเทคนิค Word2Vec
- 3) ศึกษาการแบบจำลองประเภท Sequence-to-Sequence
- 4) ศึกษาหาเครื่องมือที่เหมาะสมในการทำปัญหาพิเศษ
- 5) จัดหาข้อมูลและจัดเตรียมข้อมูลให้เหมาะสม
- 6) สร้างแบบจำลองการเรียนรู้ของเครื่องและทดสอบประสิทธิภาพ
- 7) สรุปผลและแนะนำแนวทางการพัฒนาต่อ

1.6 อุปกรณ์ที่ใช้ในการทำปัญหาพิเศษ

1) ฮาร์ดแวร์

- NVIDIA GeForce GTX 1060 หน่วยความจำ 6 GB

2) ซอฟต์แวร์

- ระบบปฏิบัติการ Windows 10 Home 64-bit
- Python 3.6.5
- TensorFlow-GPU 1.13.1
- PyThaiNLP 2.0.5
- NLTK 3.3
- CUDA Toolkit 10.0



บทที่ 2

ทฤษฎีและเอกสารที่เกี่ยวข้อง

2.1 การตัดคำ

การตัดคำ (Word Segmentation) คือกระบวนการหาขอบเขตของคำ ซึ่งวิธีการตัดคำในแต่ละภาษามีขั้นตอนวิธีต่างกันออกไป ขึ้นอยู่กับโครงสร้างและลักษณะเฉพาะของแต่ละภาษา ตัวอย่างเช่น ภาษาอังกฤษมีการเว้นระหว่างคำทำให้ง่ายต่อการตัดคำ สำหรับภาษาไทยเนื่องจากมีลักษณะการเขียนคำที่ติดต่อกันไปจนจบประโยคหรือวลีจึงทำให้ยากต่อการตัดคำ ดังนั้นการตัดคำภาษาไทย (Thai Word Segmentation) จึงมีการพัฒนาขั้นตอนวิธีต่าง ๆ เพื่อสร้างตัวตัดคำภาษาไทย (Thai Tokenizer) โดยมีความแตกต่างกันออกไปตามกาลเวลาและเพื่อให้เหมาะสมต่อการใช้งาน ซึ่งเกณฑ์การตัดคำภาษาไทย มีดังนี้

ตารางที่ 2.1 เกณฑ์การตัดคำภาษาไทย

เกณฑ์	การตัดคำ
อิงตามคำมูล	ฉั นนำ ดอก ไม้ ไป ไหว้ ศาล พระ ภูมิ ที่ โรง เรียน ประจำ
อิงตามคำมูลที่รวมกันแน่น	ฉั นนำ ดอก ไม้ ไป ไหว้ ศาล พระ ภูมิ ที่ โรง เรียน ประจำ
อิงตามคำประสมเปลี่ยนความหมาย	ฉั นนำ ดอก ไม้ ไป ไหว้ ศาล พระ ภูมิ ที่ โรง เรียน ประจำ
อิงตามหน่วยความหมาย	ฉั นนำ ดอก ไม้ ไป ไหว้ ศาล พระ ภูมิ ที่ โรง เรียน ประจำ

จากตารางที่ 2.1 แสดงเกณฑ์การตัดคำภาษาไทยและผลที่ได้จากการตัดคำตามเกณฑ์นั้น ๆ ของประโยค “ฉั|นนำ|ดอก|ไม้|ไป|ไหว้|ศาล|พระ|ภูมิ|ที่|โรง|เรียน|ประจำ” ซึ่งการพิจารณาในการใช้เกณฑ์ขึ้นอยู่กับความเหมาะสมต่อการใช้งาน ตัวอย่างเช่น การประยุกต์การตัดคำกับการแปลภาษาจะใช้เกณฑ์การตัดคำอิงตามหน่วยความหมาย เป็นต้น

2.1.1 วิธีการตัดคำ

ขั้นตอนวิธีสำหรับตัดคำมีหลากหลายวิธี ซึ่งในแต่ละขั้นตอนวิธีจะให้ประสิทธิภาพที่แตกต่างกันออกไป โดยส่วนใหญ่แล้วจะแบ่งวิธีการตัดคำภาษาไทยออกเป็น 3 ประเภท ได้แก่ ใช้กฎ (Rule-based) ใช้พจนานุกรม (Dictionary-based) และใช้คลังข้อความ (Corpus-based)

2.2.1.1 ใช้กฎ

วิธีการตัดคำโดยใช้กฎ (Rule-based) คือขั้นตอนวิธีที่พัฒนาขึ้นโดยการอาศัยหลักไวยากรณ์ภาษาไทยเพื่อสร้างกฎขึ้น ซึ่งถูกพัฒนาขึ้นในยุคแรกเริ่มของการพัฒนาตัวตัดคำภาษาไทย โดยงานวิจัยเน้นไปที่การวิเคราะห์ในระดัคำหรือพยางค์เป็นส่วนใหญ่

2.2.1.2 ใช้พจนานุกรม

วิธีการตัดคำโดยใช้พจนานุกรม (Dictionary-based) คือขั้นตอนวิธีที่พัฒนาขึ้นโดยใช้สายอักขระ (String) เปรียบเทียบกับคำที่จัดเก็บในพจนานุกรม

2.2.1.3 ใช้คลังข้อความ

วิธีการตัดคำโดยใช้คลังข้อความ (Corpus-based) คือขั้นตอนวิธีที่พัฒนาขึ้นโดยการสร้างแบบจำลองภาษา (Language Model) ด้วยเทคนิคต่าง ๆ ด้วยคลังข้อความ เช่น วิธีการทางสถิติ การเรียนรู้ของเครื่องจักร เป็นต้น

ตารางที่ 2.2 การเปรียบเทียบข้อดีและข้อเสียของวิธีการตัดคำภาษาไทย

วิธีการ	ข้อดี	ข้อเสีย
ใช้กฎ	1) สามารถตัดคำระดับพยางค์ได้ดี	1) กฎสามารถสร้างได้ยาก
ใช้พจนานุกรม	1. สามารถเพิ่มคำศัพท์ใหม่ในพจนานุกรมได้	1) การรวบรวมพจนานุกรมต้องใช้ทรัพยากรมนุษย์ในการจัดการสูง
ใช้คลังข้อความ	1) สามารถสร้างกฎต่าง ๆ ได้โดยไม่ต้องวิเคราะห์หลักภาษา 2) สามารถตัดคำศัพท์ที่ไม่รู้จักหรือคำผิดได้บางกรณี	1) ระยะเวลาในการเรียนรู้สูงเนื่องจากความซับซ้อนของขั้นตอนวิธี 2) ต้องการคลังข้อความที่ถูกต้องและเหมาะสมจำนวนมาก

2.1.2 เทคนิคการตัดคำ

เทคนิคสำหรับการตัดคำในภาษาไทยถูกพัฒนาอย่างต่อเนื่อง แต่ละเทคนิคให้ผลลัพธ์ที่แตกต่างกันดังตัวอย่างการตัดคำภาษาไทยในตารางที่ 2.1 โดยในภาววิจยนี้มีเทคนิคที่เกี่ยวข้องกับการตัดคำภาษาไทย ดังนี้

2.1.2.1 เทคนิคการตัดคำด้วยกลุ่มตัวอักษรไทย [1]

เทคนิคการตัดคำด้วยกลุ่มตัวอักษรไทย (Thai Character Cluster: TCC) [1] เป็นขั้นตอนวิธีที่ใช้การแยกแต่ละอักขระเพื่อทำการวิเคราะห์ก่อนทำการกำหนดตำแหน่งขอบเขตของคำ ซึ่งอาศัยแนวคิดว่ามีตัวอักษรบางตัวที่ทราบว่าจะไม่ควรตัดคำตำแหน่งนั้น ตัวอย่างเช่น C1 จะไม่สามารถตัดคำหน้าสระอาได้ ซึ่งสามารถแก้ปัญหาคำที่ไม่รู้จัก (Unknown Word) ได้ในบางกรณี

2.1.2.2 เทคนิคการตัดคำแบบเลือกคำยาวที่สุด

เทคนิคการตัดคำแบบเลือกคำยาวที่สุด (Longest Matching) เป็นเทคนิคที่ใช้หลักการตัดคำโดยใช้พจนานุกรม โดยใช้สายอักขระ (String) เปรียบเทียบกับคำที่จัดเก็บในพจนานุกรม โดยเลือกคำที่มีความยาวมากที่สุดในพจนานุกรม ตัวอย่างเช่น

ตาม|หาม|เหสี = ตาม|หาม|เห|สี|

เอนั่ง|ตาก|ลม|ที่|ชาย|หาด = เอนั่ง|ตาก|ลม|ที่|ชาย|หาด|

2.1.2.3 เทคนิคการตัดคำแบบสอดคล้องมากที่สุด

เทคนิคการตัดคำแบบสอดคล้องมากที่สุด (Maximal Matching) เป็นเทคนิคที่พัฒนาขึ้นเพื่อแก้ปัญหาของเทคนิคการตัดคำแบบเลือกคำยาวที่สุด ซึ่งมีข้อผิดพลาด

จากการเลือกคำที่มีความยาวมากที่สุด เช่น ตาม|หาม|เห|สี|

เทคนิคการตัดคำแบบสอดคล้องมากที่สุดจะทำการค้นหาเส้นทางของการแบ่งคำที่เป็นไปได้ทั้งหมด และทำการเลือกเส้นทางที่มีจำนวนคำที่น้อยที่สุด ตัวอย่างเช่น

$$\begin{aligned} \text{ตาม|หาม|เหสี} &= \text{ตาม|หาม|เหสี} = 4 \text{ คำ} \\ &= \text{ตาม|หาม|เหสี} = 3 \text{ คำ (ถูกเลือก)} \end{aligned}$$

ในกรณีที่เส้นทางการแบ่งคำทั้งหมดมีจำนวนของคำเท่ากันจะใช้เทคนิคการตัดคำแบบยาวที่สุด ตัวอย่างเช่น

$$\begin{aligned} \text{เธอนั่งตากลมที่ชายหาด} &= \text{เธอนั่ง|ตาก|ลม|ที่|ชายหาด} = 7 \text{ คำ (ถูกเลือก)} \\ &= \text{เธอนั่ง|ตาก|ลม|ที่|ชายหาด} = 7 \text{ คำ} \end{aligned}$$

2.2 Word2Vec [2]

Word2Vec คือการนำเสนอคำ (Word Represent) ในรูปแบบเวกเตอร์ โดยจะมีคุณสมบัติเหมือนกับค่าเวกเตอร์ เช่น ค่าของเวกเตอร์มีความใกล้เคียงกันมาก หมายความว่าคำนั้น ๆ มีความสัมพันธ์กันมาก

การสร้างค่าเวกเตอร์ของคำสามารถทำได้โดยการใช้โครงข่ายประสาทเทียมร่วมกับแบบจำลองทางภาษา (Language Model) ซึ่งมี 2 แบบจำลอง คือ Skip-Gram model และ CBOW (Continuous Bag-of-Words) ซึ่งมีจุดประสงค์ในการใช้งานแบบจำลองที่แตกต่างกันดังนี้

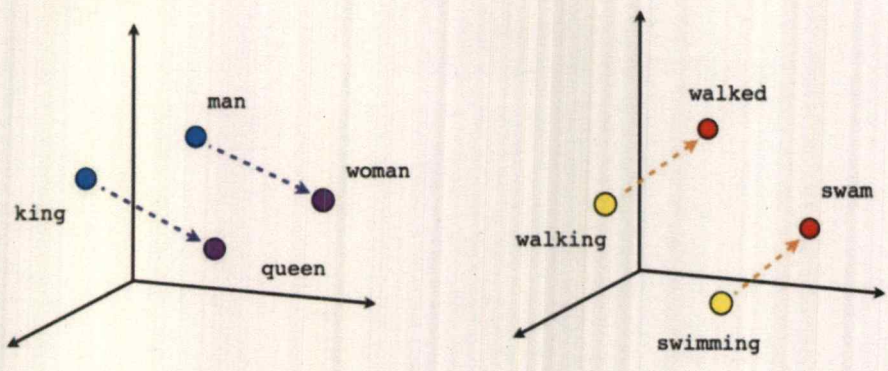
การสร้างค่าเวกเตอร์ด้วยแบบจำลอง CBOW ในการประยุกต์ใช้งานจะใช้เพื่อการทำนายคำที่ขาดหายไปบริบท ตัวอย่างเช่น

ประโยคที่ 1: อยากไปเที่ยว ___ ที่มีริมชายหาดสวยๆ

ประโยคที่ 2: อากาศร้อนมากอยาก ___ อะไรที่มันเย็นๆ

จากประโยคที่ 1 และประโยคที่ 2 ในช่องว่างมีคำที่ขาดหายไปซึ่งเมื่อพิจารณาจากบริบทแล้วสามารถมีคำที่นำมาเติมได้หลายคำ ซึ่งการสร้างแบบจำลอง Word2Vec ด้วยแบบจำลอง CBOW จะสามารถค้นหาคำที่เหมาะสมต่อบริบทรอบข้างได้

การสร้างค่าเวกเตอร์ด้วยแบบจำลอง Skip-Gram ในการประยุกต์ใช้งานจะใช้เพื่อการทำนายบริบทของคำว่ามีความเกี่ยวข้องกับบริบทใด ซึ่งบริบทคือคำหรือข้อความแวดล้อมเพื่อช่วยให้เข้าใจความหมาย กล่าวคือ การประยุกต์ใช้งานจะใช้เพื่อการทำนายว่าคำหนึ่งๆมีความเกี่ยวข้องกับคำใดบ้าง และจากหลักสำคัญนี้เองทำให้สามารถทราบได้ว่าคำใด ๆ ในภาษามีความสัมพันธ์กันอย่างไร เช่น มีความคล้ายคลึงกันมากน้อยเพียงใด และด้วยความสามารถของแบบจำลองนี้ ทำให้ได้รับความนิยมและนำประยุกต์ใช้งานอย่างแพร่หลาย เช่น การจำแนกประเภทข้อความ (Text Classification) เป็นต้น



Male-Female

Verb tense

รูปที่ 2.1 แสดงความสัมพันธ์ของคำ

(ที่มา: <https://towardsdatascience.com/word-embedding-with-word2vec-and-fasttext-a209c1d3e12c>)

จากรูปที่ 2.1 ด้านซ้ายแสดงความสัมพันธ์ของคำว่า “king”, “man”, “queen”, “woman” โดยแสดงความสัมพันธ์ของคำว่า “king”, “man” เป็นเพศชาย “queen”, “woman” เป็นเพศหญิง หรือ “king”, “queen” จะใช้ในบริบทเดียวกันและ “man”, “woman” จะใช้ในบริบทเดียวกัน

จากรูปที่ 2.1 ด้านขวาแสดงความสัมพันธ์ของคำว่า “walking”, “walked”, “swimming”, “swam” โดยแสดงความสัมพันธ์ของคำว่า “walking”, “walked” มีรากศัพท์เป็นคำเดียวกัน “swimming”, “swam” มีรากศัพท์เป็นคำเดียวกัน หรือ “walking”, “swimming” ใช้ในรูปแบบของกริยา (Tense) เดียวกัน และ “walked”, “swam” ใช้ในรูปแบบของกริยาเดียวกัน

2.3 ขนาดของแบตช์ (Batch Size)

คือการแบ่งข้อมูลออกเป็นส่วนๆ เพื่อการจำกัดการใช้ทรัพยากรการประมวลผล เนื่องจากจำนวนข้อมูลในการสร้างแบบจำลองมีจำนวนมาก แต่ในการสร้างแบบจำลองมีทรัพยากรไม่เพียงพอสำหรับนำข้อมูลเข้าสู่หน่วยคำจำเพื่อประมวลผล ดังนั้นจึงต้องแบ่งข้อมูลออกเป็นส่วนๆ เช่น ขนาดของแบตช์เท่ากับ 1,000 แปลว่า นำข้อมูลเข้าสู่หน่วยคำจำครั้งละ 1,000 แถวและแบบจำลองมีการเรียนรู้ข้อมูลครั้งละ 1,000 ข้อมูล

ทั้งนี้ในการทำงานครบ 1 แบตช์ยังนิยมให้มีการคำนวณข้อผิดพลาดของแบบจำลอง (Loss Value) และปรับพารามิเตอร์ (Parameter) ของแบบจำลองให้มีความถูกต้องมากยิ่งขึ้น ซึ่งความถี่ในการปรับค่าพารามิเตอร์จะส่งผลต่อความถูกต้องของแบบจำลอง ดังนั้นแล้วจึงต้องมีการกำหนดค่าความถี่ในการปรับค่าพารามิเตอร์ให้เหมาะสม

2.4 Epoch

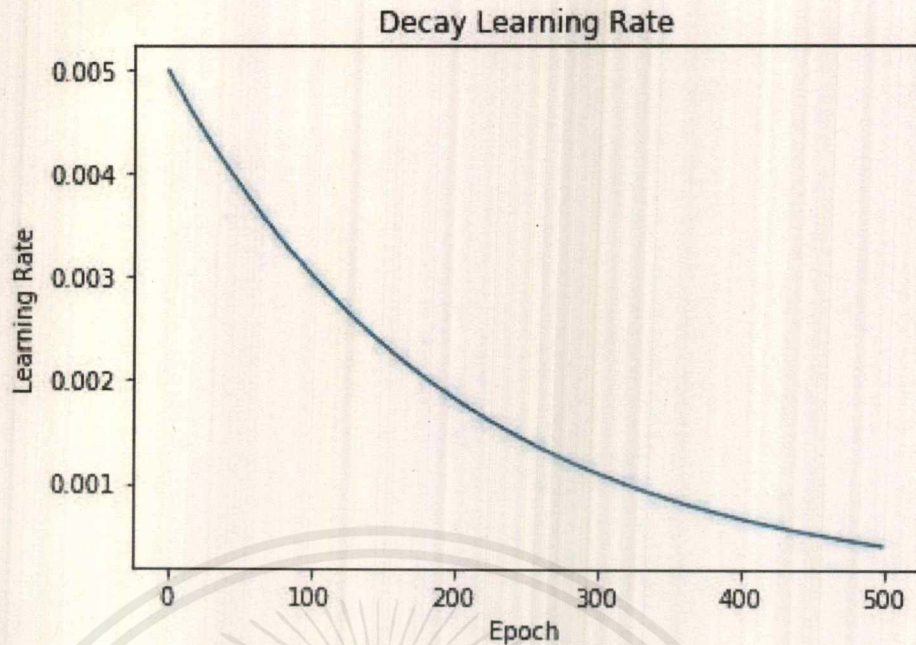
คือ 1 รอบในการเรียนรู้ของแบบจำลองกับข้อมูลทั้งหมด เช่น จำนวน Epoch เท่ากับ 5 คือ กำหนดให้แบบจำลองมีการเรียนรู้กับข้อมูลทั้งหมดจำนวน 5 รอบ

2.5 อัตราการเรียนรู้ (Learning Rate)

คือ ค่าที่บ่งบอกถึงความรวดเร็วในการปรับตัวตามข้อมูลการสอนของแบบจำลอง หากมีอัตราการเรียนรู้น้อยเกินไปจะทำให้แบบจำลองเข้าใกล้คำตอบที่ถูกต้องช้า แต่หากมีอัตราการเรียนรู้ที่มากเกินไปจะทำให้แบบจำลองจำเพาะเจาะจงกับข้อมูลฝึกสอนล่าสุดมากเกินไป ดังนั้นแล้วต้องมีการกำหนดค่าอัตราการเรียนรู้ให้มีความเหมาะสมกับแต่ละแบบจำลอง

2.6 การลดอัตราการเรียนรู้ (Decay Learning Rate)

จากปัญหาในการปรับค่าความเหมาะสมของอัตราการเรียนรู้ การลดอัตราการเรียนรู้จะสามารถลดปัญหาการกำหนดค่าอัตราการเรียนรู้ที่เหมาะสมได้ โดยนิยมกำหนดให้มีการลดอัตราการเรียนรู้คือ เมื่อเริ่มต้นจะกำหนดให้อัตราการเรียนรู้มีค่ามากเพื่อให้แบบจำลองมีการเรียนรู้ที่รวดเร็ว จากนั้นทำการปรับลดค่าอัตราการเรียนรู้ลงเรื่อย ๆ เพื่อให้แบบจำลองมีการเรียนรู้ที่ช้าลงเพื่อเข้าใกล้คำตอบที่ถูกต้องมากที่สุด



รูปที่ 2.2 แสดงการลดลงของอัตราการเรียนรู้

2.7 One-Hot Encoding

คือชุดของบิต (1 หรือ 0) ที่นำเสนอแทนข้อมูลต่าง ๆ โดยมีค่า 1 เพียง 1 ตำแหน่งและตำแหน่งอื่น ๆ มีค่าเท่ากับ 0 ตัวอย่างเช่น “ฉันรักเมืองไทย” สามารถแปลงค่าเป็น One-Hot ของแต่ละคำได้ดังนี้

ฉัน	=	[1, 0, 0, 0]
รัก	=	[0, 1, 0, 0]
เมือง	=	[0, 0, 1, 0]
ไทย	=	[0, 0, 0, 1]

2.8 ฟังก์ชันความสูญเสีย (Loss Function)

ฟังก์ชันความสูญเสีย คือฟังก์ชันที่ใช้วัดความแตกต่างระหว่างค่าจริงกับค่าที่แบบจำลองทำนาย

2.9 Categorical Cross-Entropy [3]

คือฟังก์ชันที่ใช้ค่าครอสเอนโทรปี (Cross-Entropy) ในการวัดความสูญเสียจากการทำนายของแบบจำลอง

เอนโทรปี คือการวัดความไม่แน่นอนในการทำนายผลลัพธ์ของแบบจำลอง

$$H(p) = -\sum_{i=0}^N p_i \log p_i \quad (2.1)$$

เมื่อ p คือความน่าจะเป็นที่จะเกิดคลาสต่าง ๆ จากทั้งหมด N คลาส

H คือค่า Entropy ของ p

Cross-Entropy คือการวัดค่าความต่างของการแจกแจงความน่าจะเป็นในการทำนายผลลัพธ์ของแบบจำลอง

$$H(p, q) = -\sum_{i=0}^N p_i \log q_i \quad (2.2)$$

เมื่อ p คือความน่าจะเป็นของผลลัพธ์ในแต่ละคลาส (Class) จากคลาสทั้งหมด N คลาส

q คือผลลัพธ์ในการทำนายคลาส (Class) ซึ่งเป็นรูปแบบของ One-Hot Encoding

H คือค่า Cross-Entropy ของ p และ q

2.10 Optimizer

Optimizer คืออัลกอริทึมสำหรับปรับค่าถ่วงน้ำหนัก (Weight) เพื่อให้แบบจำลองการเรียนรู้ด้วยเครื่องมีค่าสูญเสีย (Loss) ความถูกต้องน้อยที่สุด

2.11 ROUGE (Recall-Oriented Understudy for Gisting Evaluation) [4]

คือตัวชี้วัดสำหรับประเมินผลการสรุปข้อความอัตโนมัติ (Automatic Summarization) และการแปลภาษาด้วยเครื่อง (Machine Translation) โดยการอ้างอิงจากการสรุปหรือการแปลที่มนุษย์สร้างขึ้น โดยมีสมการดังนี้

$$Recall = \frac{\text{จำนวนคำที่คาบเกี่ยวกัน (Overlap)}}{\text{จำนวนคำในประโยคสรุปข่าวที่ถูกต้อง}} \quad (2.3)$$

ค่า Recall คือ การวัดความสามารถในการดึงข้อมูลที่เกี่ยวข้องออกมา หรือสร้างข้อความผลลัพธ์ที่ถูกต้อง

$$Precision = \frac{\text{จำนวนคำที่คาบเกี่ยวกัน (Overlap)}}{\text{จำนวนคำในประโยคสรุปข่าวที่แบบจำลองสร้างขึ้น}} \quad (2.4)$$

ค่า Precision คือการวัดความสามารถในการขจัดข้อมูลที่ไม่เกี่ยวข้องออกไป หรือไม่สร้างข้อความผลลัพธ์ที่ไม่ถูกต้อง

2.11.1 ROUGE-1

คือการวัดความถูกต้องของการสรุปข้อความหรือการแปลภาษาโดยวัดจากความคาบเกี่ยวกัน (Overlap) ของคำด้วยความยาว 1 คำ ตัวอย่างเช่น

ประโยคสรุปข่าวที่แบบจำลองสร้างขึ้น: the cat was found under the bed

ประโยคสรุปข่าวที่ถูกต้อง: the cat was under the bed

มีค่า Recall เท่ากับ

$$Recall = \frac{6}{6} = 1.0$$

มีค่า Precision เท่ากับ

$$Precision = \frac{6}{7} = 0.86$$

2.11.2 ROUGE-2

คือการวัดความถูกต้องของการสรุปข้อความหรือการแปลภาษาโดยวัดจากความคาบเกี่ยวกัน (Overlap) ของคำด้วยความยาวจำนวน 2 คำ จากประโยคสรุปข่าวที่แบบจำลองสร้างขึ้นและประโยคสรุปข่าวที่ถูกต้อง จะได้ลากับคู่ดังนี้

the	cat	the	cat
cat	was	cat	was
was	found	was	under
found	under	under	the
under	the	the	bed
the	bed		

มีค่า Recall เท่ากับ

$$ROUGE2_{Recall} = \frac{4}{5} = 0.8$$

มีค่า Precision เท่ากับ

$$ROUGE2_{Precision} = \frac{4}{6} = 0.67$$

2.11.3 ROUGE-L

คือการวัดความถูกต้องของการสรุปข้อความหรือการแปลภาษา ซึ่งมีความคล้ายคลึงกับ ROUGE-1 และ ROUGE-2 ดังที่ได้กล่าวมาในข้างต้น แต่สำหรับ ROUGE-L จะใช้สมการ LCS (Longest Common Subsequence) ในการคำนวณ ซึ่งข้อดีของการใช้ LCS คือไม่ต้องใช้การจับคู่แบบต่อเนื่อง แต่เป็นการจับคู่เรียงลำดับที่สะท้อนถึงการเรียงลำดับคำในระดับประโยค ซึ่งมีสมการดังนี้

$$Recall = \frac{LCS(\text{ประโยคแบบจำลองสร้างขึ้น, ประโยคเป้าหมาย})}{\text{จำนวนคำในประโยคสรุปข่าวที่ถูกต้อง}} \tag{2.5}$$

$$Precision = \frac{LCS(\text{ประโยคแบบจำลองสร้างขึ้น, ประโยคเป้าหมาย})}{\text{จำนวนคำในประโยคสรุปข่าวที่แบบจำลองสร้างขึ้น}} \tag{2.6}$$

2.12 F-measure

F-measure (F_1 Score หรือ F Score) คือค่าเฉลี่ยประสิทธิโดยรวมของแบบจำลองการเรียนรู้ของเครื่อง

$$F = \frac{2 \times Precision \times Recall}{Precision + Recall} \tag{2.7}$$

2.13 Longest Common Subsequence Problem (LCS) [5]

เมื่อต้องการเปรียบเทียบชุดอักษร 2 ชุด ได้แก่ “AGCAT” และ “GAC” ด้วย LCS จะมีขั้นตอนดำเนินการ ดังนี้

ตารางที่ 2.3 เริ่มต้นสร้างตาราง LCS

	∅	A	G	C	A	T
∅	∅	∅	∅	∅	∅	∅
G	∅					
A	∅					
C	∅					

1. เริ่มต้นนำค่าที่ต้องการคำนวณมาเขียนเป็นตารางดังเช่นตารางที่ 2.3
2. เริ่มเปรียบเทียบอักษรในแถวที่ 2 คือ G คือ กับอักษรคอลัมน์ที่ 2 คือ A ว่าอักษรเหมือนกันหรือไม่ ซึ่งได้ผลลัพธ์คือ ไม่เหมือนกัน และนำ ∅ พร้อมทั้งลูกศรชี้ขึ้นและลงเพิ่มในตารางดังรูปอักษร
3. จากนั้นดำเนินการต่อกับคอลัมน์ที่ 3 จะได้ผลลัพธ์คืออักษรเหมือนกัน ให้ทำการเพิ่มอักษร G และลูกศรชี้ขึ้นเฉียงไปทางด้านซ้ายมือในตาราง
4. จากนั้นดำเนินการต่อกับคอลัมน์ที่ 4 จะได้ผลลัพธ์คืออักษรไม่เหมือนกัน ให้ทำการเพิ่มชุดอักษรตัวก่อนหน้าที่เหมือนกัน คือ G และลูกศรเข้าหาชุดอักษรที่เหมือนกันก่อนหน้าคือชี้ไปทางซ้ายเนื่องจากชุดอักษรที่เหมือนกันก่อนหน้าอยู่ทางด้านซ้าย
5. ดำเนินการเช่นไปเรื่อย ๆ จนจบทุกตัวอักษร จะได้ผลลัพธ์ดังนี้ตารางที่ 2.4

ตารางที่ 2.4 การคำนวณ LCS กับข้อมูลแถวที่ 2

	∅	A	G	C	A	T
∅	∅	∅	∅	∅	∅	∅
G	∅	$\begin{matrix} \uparrow \\ \leftarrow \end{matrix} \emptyset$	$\swarrow (G)$	$\leftarrow (G)$	$\leftarrow (G)$	$\leftarrow (G)$
A	∅					
C	∅					

6. ดำเนินการกับแถวที่ 3 คืออักษร A กับคอลัมน์ที่ 2 คือ A จะได้ผลลัพธ์คืออักษรเหมือนกัน ให้ทำการเพิ่มอักษร A และลูกศรชี้ขึ้นเอียงไปทางด้านซ้ายมือในตาราง
7. จากนั้นดำเนินการต่อกับคอลัมน์ที่ 3 จะได้ผลลัพธ์คืออักษรไม่เหมือนกัน ให้ทำการเพิ่มชุดอักษรตัวก่อนหน้าเหมือนกัน คือ A และ G ซึ่งจะลูกศรชี้เข้าหาชุดอักษรที่เหมือนกันก่อนหน้า คือทางด้านซ้ายเพื่อชี้เข้าหา A และทางด้านบน เพื่อชี้เข้าหา G และสำหรับการดำเนินการกับคอลัมน์ที่ 4 จะได้ผลลัพธ์เช่นกับกับการดำเนินการกับคอลัมน์ที่ 3
8. ในการดำเนินการกับคอลัมน์ที่ 5 จะได้ผลลัพธ์คือมีอักษรเหมือนกัน ให้ทำการเพิ่มอักษร GA และลูกศรชี้ขึ้นเอียงไปทางด้านซ้ายมือในตาราง เนื่องจากจากก่อนหน้านี้มีรูปแบบอักษรเหมือนกันคือ G และในการดำเนินการนี้มีอักษรที่เหมือนกันคือ A ดังนั้น ชุดอักษรที่ต้องเพิ่มในตารางคือ GA
9. ดำเนินการเช่นนี้ไปเรื่อย ๆ จนจบทุกตัวอักษร จะได้ผลลัพธ์ดังนี้ตารางที่ 2.5

ตารางที่ 2.5 การคำนวณ LCS กับข้อมูลแถวที่ 3

	\emptyset	A	G	C	A	T
\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset
G	\emptyset	$\begin{matrix} \uparrow \\ \leftarrow \emptyset \end{matrix}$	$\swarrow (G)$	$\leftarrow (G)$	$\leftarrow (G)$	$\leftarrow (G)$
A	\emptyset	$\swarrow (A)$	$\begin{matrix} \uparrow \\ \leftarrow (A) \end{matrix} \& (G)$	$\begin{matrix} \uparrow \\ \leftarrow (A) \end{matrix} \& (G)$	$\swarrow (GA)$	$\leftarrow (GA)$
C	\emptyset					

10. เมื่อดำเนินการครบทุกตัวอักษรแล้ว จะได้ผลลัพธ์ดังนี้

ตารางที่ 2.6 การคำนวณ LCS กับข้อมูลแถวที่ 4

	\emptyset	A	G	C	A	T
\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset
G	\emptyset	$\begin{matrix} \uparrow \\ \leftarrow \emptyset \end{matrix}$	$\swarrow (G)$	$\leftarrow (G)$	$\leftarrow (G)$	$\leftarrow (G)$
A	\emptyset	$\swarrow (A)$	$\begin{matrix} \uparrow \\ \leftarrow (A) \end{matrix} \& (G)$	$\begin{matrix} \uparrow \\ \leftarrow (A) \end{matrix} \& (G)$	$\swarrow (GA)$	$\leftarrow (GA)$
C	\emptyset	$\begin{matrix} \uparrow \\ \leftarrow (A) \end{matrix}$	$\begin{matrix} \uparrow \\ \leftarrow (A) \end{matrix} \& (G)$	$\swarrow (AC) \& (GC)$	$\begin{matrix} \uparrow \\ \leftarrow (AC) \end{matrix} \& (GC) \& (GA)$	$\begin{matrix} \uparrow \\ \leftarrow (AC) \end{matrix} \& (GC) \& (GA)$

11. จากนั้นทำการนับจำนวนความยาวสูงสุดของชุดอักษรในแถวแต่ละช่องตาราง จะได้ผลลัพธ์ดังเช่นในตาราง 2.7

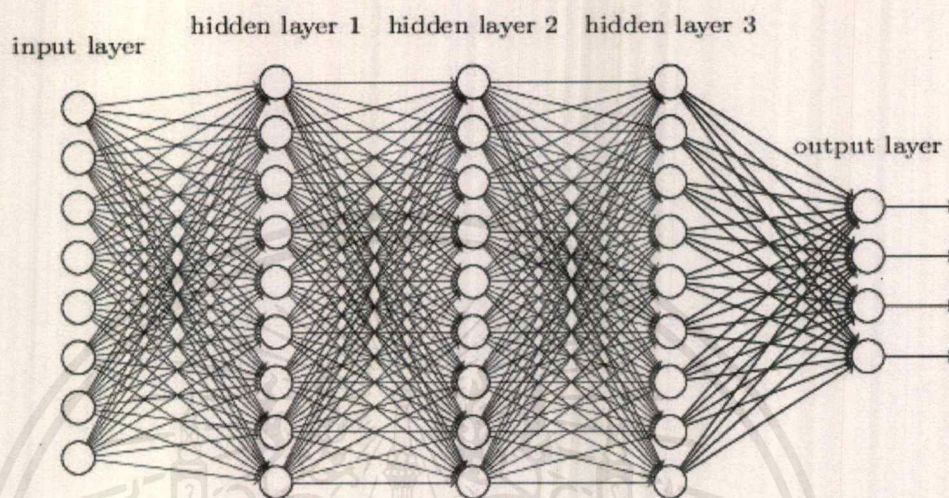
ตารางที่ 2.7 จำนวนความยาวสูงสุดของชุดอักขรในแต่ละช่องตาราง

	∅	A	G	C	A	T
∅	0	0	0	0	0	0
G	0	$\begin{matrix} \uparrow \\ \leftarrow 0 \end{matrix}$	$\swarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$
A	0	$\swarrow 1$	$\begin{matrix} \uparrow \\ \leftarrow 1 \end{matrix}$	$\begin{matrix} \uparrow \\ \leftarrow 1 \end{matrix}$	$\swarrow 2$	$\leftarrow 2$
C	0	$\uparrow 1$	$\begin{matrix} \uparrow \\ \leftarrow 1 \end{matrix}$	$\swarrow 2$	$\begin{matrix} \uparrow \\ \leftarrow 2 \end{matrix}$	$\begin{matrix} \uparrow \\ \leftarrow 2 \end{matrix}$

จากตารางจะได้ผลลัพธ์การดำเนินการของฟังก์ชัน LCS เท่ากับ 2 (ค่าที่ด้านล่างมุมขวาของตาราง) ซึ่งมีความหมายว่า มีชุดอักขรที่เหมือนกันยาวที่สุดเท่ากับ 2 ตัวอักษร ได้แก่ AC, GC และ GA

2.14 การเรียนรู้เชิงลึก (Deep Learning) [6]

การเรียนรู้เชิงลึกเป็นสาขาของการเรียนรู้ของเครื่อง หลักสำคัญของการเรียนรู้เชิงลึก คือ อัลกอริทึมที่พยายามสร้างแบบจำลองเพื่อแทนรูปแบบและความสัมพันธ์ของข้อมูลที่ซับซ้อน โดยการสร้างสถาปัตยกรรมข้อมูลและประกอบไปด้วยโครงสร้างย่อย ๆ หลายอันหรือหลายชั้นเชื่อมต่อกัน



รูปที่ 2.3 โครงข่ายประสาทเทียมแบบลึก

(ที่มา: <https://www.kdnuggets.com/2017/05/deep-learning-big-deal.html>)

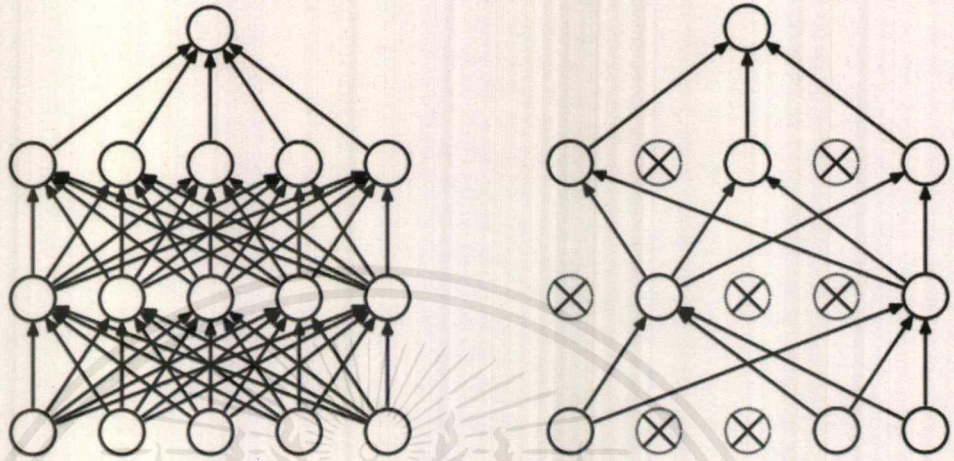
จากรูป คือ โครงข่ายประสาทเทียมแบบลึก (Deep Artificial Neural Networks) ประกอบไปด้วยชั้นอินพุต, ชั้นเอาต์พุตและชั้นซ่อนจำนวน 2 ชั้น ซึ่งเป็นหนึ่งในสถาปัตยกรรมของการเรียนรู้เชิงลึก (Deep Learning)

2.15 การเข้ากันมากเกินไป (Overfitting)

Overfitting คือ แบบจำลองสามารถจดจำข้อมูลชุดฝึกสอนได้เป็นอย่างดี แต่เมื่อนำแบบจำลองไปทดสอบกับข้อมูลชุดทดสอบแล้วจะได้ผลลัพธ์ที่มีประสิทธิภาพต่ำ เนื่องจากแบบจำลองสามารถเข้ากัน (Fitting) กับข้อมูลชุดฝึกสอนได้เป็นอย่างดีแต่ไม่เข้ากันหรือไม่สามารถคาดการณ์ (Prediction) ผลลัพธ์ที่ถูกต้องของข้อมูลชุดทดสอบได้

2.16 Dropout [7]

คือการปิดบังนิวรอนของโครงข่ายประสาทเทียมในบางส่วน เพื่อเพิ่มความยากในการเรียนรู้ของแบบจำลอง ซึ่งวิธีการ Dropout นี้ใช้เพื่อลดปัญหา Overfitting ดังที่ได้กล่าวมาในข้างต้น



รูปที่ 2.4 แสดงการ Dropout ของแบบจำลองโครงข่ายประสาทเทียม

(ที่มา: <http://jmlr.org/papers/volume15/srivastava14a/srivastava14a.pdf>)

จากรูปที่ 2.4 ทางด้านซ้ายคือแบบจำลองโครงข่ายประสาทเทียมโดยปกติที่ไม่มี การ Dropout และรูปทางด้านขวาคือ แบบจำลองโครงข่ายประสาทเทียมที่มีการ Dropout โดยจำนวนของการปิดบังนิวรอนขึ้นอยู่กับข้อกำหนดค่าเปอร์เซ็นต์การ Dropout ของผู้พัฒนา ซึ่งจะส่งผลต่อประสิทธิภาพการทำงานของแบบจำลองการเรียนรู้ของเครื่อง

2.17 การเรียนแบบวนซ้ำ (Recurrent Neural Network: RNN) [8]

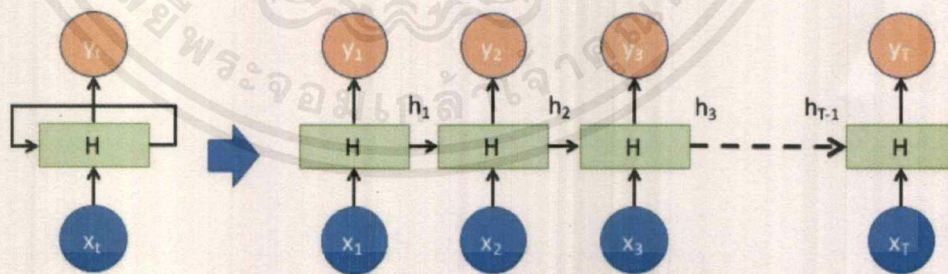
สถาปัตยกรรมของ RNN ถูกออกแบบมาเพื่อใช้งานกับข้อมูลที่มีลักษณะเป็นลำดับ (Sequential Data) เช่น ข้อความ (Sequence of Words) เมื่อพิจารณาพฤติกรรมการอ่านหนังสือของคนเราจะได้ว่า เราจะอ่านทีละคำจากซ้ายไปขวา (สำหรับภาษาไทยหรือภาษาอังกฤษ) และการที่เราสามารถเข้าใจได้ว่าประโยคเหล่านั้นมีความเกี่ยวข้องกับเรื่องใด เนื่องจากเราได้นำเอาเรื่องราวจากสิ่งที่ได้อ่านผ่านไปแล้ว (สถานะซ่อนหรือสถานะก่อนหน้า) มาผสมกับคำที่เพิ่งอ่านไป (ข้อมูลอินพุตหรือคำที่เรา กำลังอ่าน ณ เวลานั้น) ทำให้สามารถเข้าใจความหมายในส่วนตรงที่กำลังอ่านได้ ซึ่งในสถาปัตยกรรม RNN จะใช้หลักการเดียวกันกับการอ่านหนังสือนี้ คือการปรับรูปแบบของโครงข่ายประสาทเทียมเดิมเพื่อให้สามารถนำเอาสถานะ (State) หรือความรู้ก่อนหน้ามารวมกับข้อมูลอินพุตตัวใหม่ที่เพื่อทำความเข้าใจข้อมูลใหม่ไปเรื่อย ๆ

ดังนั้นสถาปัตยกรรม RNN มี 2 ส่วนที่สำคัญคือ

1. สถานะซ่อน (Hidden State) ก่อนหน้า
2. ข้อมูลอินพุต ณ เวลานั้น

สำหรับการอธิบายถึงการทำงานต่อไป กำหนดให้มีตัวแปรดังนี้

- ข้อมูลอินพุต คือ $x_1, x_2, x_3, \dots, x_t$
- สถานะซ่อนที่เวลา t จะใช้ตัวแปรว่า h_t
- ถ้าข้อมูลอินพุตคือประโยคว่า “ฉันกินข้าว” จะได้ว่า $x_1 = \text{“ฉัน”}$, $x_2 = \text{“กิน”}$ และ $x_3 = \text{“ข้าว”}$



รูปที่ 2.5 สถาปัตยกรรม RNN

(ที่มา: <https://medium.com/@sinart.t/long-short-term-memory-lstm-e6cb23b494c6>)

โดยที่

H = Hidden Layer

y_t = เอาต์พุต (Output) จาก RNN ที่เวลา t

x_t = ข้อมูลอินพุต (Input Data) ที่เวลา t

h_t = สถานะซ่อน (Hidden State) ที่เวลา t

จากรูป 2.5 ทางซ้ายแสดงให้เห็นว่ามีการวนซ้ำ (Loop) กลับเข้ามาที่ชั้นซ่อน (Hidden Layer) ของโครงข่ายประสาทเทียมดังที่ได้กล่าวไปในข้างต้นแล้วว่าสิ่งสำคัญสิ่งขอสถาปัตยกรรม RNN คือสถานะก่อนหน้าข้อมูลอินพุต ณ เวลานั้น ซึ่งหลักสำคัญของการวนซ้ำนี้เพื่อนำสถานะซ่อนก่อนหน้ากลับมาใช้นั่นเอง หรืออาจมองได้ว่า สถาปัตยกรรม RNN คือโครงข่ายประสาทเทียมที่มีหน่วยความจำ (Memory) เพิ่มขึ้นมาเพื่อเก็บค่าสถานะซ่อนที่คำนวณเอาไว้ก่อนหน้านี้

จากรูป 2.5 ทางขวาเป็นรูปที่คล้ายรูปทางซ้ายออกเพื่อแสดงการทำงานที่ละขั้นตอนของสถาปัตยกรรม RNN

$$h_t = f_h(U_h h_{t-1} + W_h x_t + b_h) \quad (2.8)$$

$$y_t = f_y(W_y h_t + b_y) \quad (2.9)$$

โดยที่

f_h คือฟังก์ชันก่กระตุ้น (Activation Function) ของชั้นซ่อน เช่น

Tanh function [9], ReLU Function [10], Sigmoid Function [11]

f_y คือฟังก์ชันก่กระตุ้นของชั้นเอาต์พุต เช่น Softmax Function [12]

W_h คือเมทริกซ์ค่าถ่วงน้ำหนัก (Weight Matrix) ของชั้นซ่อน

U_h คือเมทริกซ์ที่เชื่อมระหว่างสถานะซ่อน (Hidden-State-to-Hidden-State Matrix หรือ Transition Matrix)

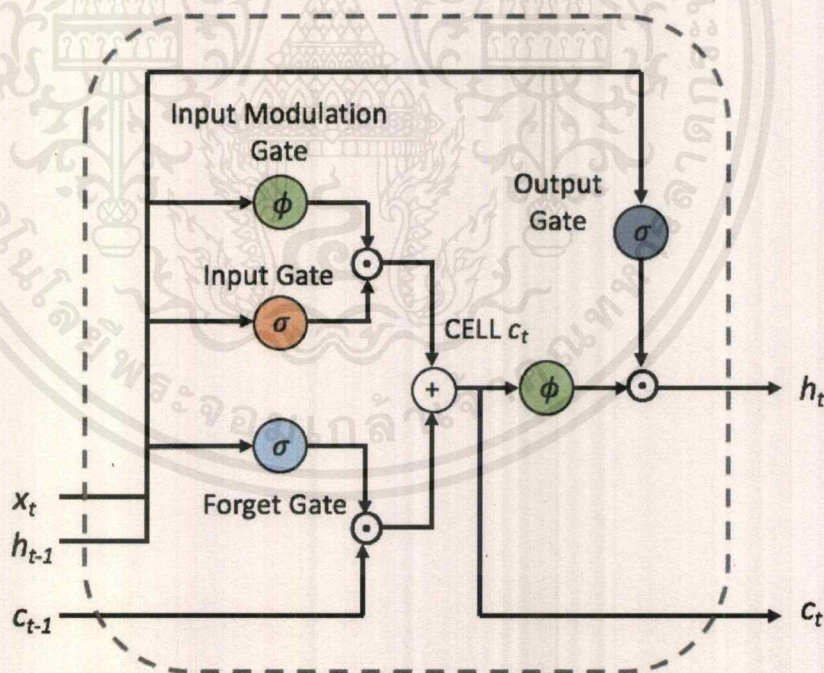
จากสมการจะเห็นว่า การคำนวณสถานะซ่อนที่เวลา t ได้นั้น (h_t) จำเป็นต้องใช้ 2 ตัวแปรสำคัญ คือสถานะซ่อนก่อนหน้าก่อนหน้า (h_{t-1}) และข้อมูลอินพุต ณ เวลานั้น (x_t)

2.17.1 ปัญหาของ RNN

ในการปรับปรุงค่าถ่วงน้ำหนักของโครงข่ายประสาทเทียมจะอาศัยอัลกอริทึม Backpropagation [13] สำหรับสถาปัตยกรรม RNN จะมีการคำนวณที่ความซับซ้อนมาก เนื่องจากเอาต์พุต y_t ไม่ได้ถูกคำนวณจากช่วงเวลา $t=t$ เพียงอย่างเดียว แต่ถูกคำนวณต่อเนื่องจาก $t=t-1, t=t-2, t=t-3, \dots, t=1$ ซึ่งจะทำให้เกิดปัญหาการสูญหายของค่าเกรเดียนต์ (Vanishing Gradient Problem) [14] และหากมีความยาวของข้อมูลมากจะทำให้มีโอกาสเกิดปัญหาการสูญหายของค่าเกรเดียนต์มาก โดยจะส่งผลให้แบบจำลองไม่เกิดการเรียนรู้ เนื่องจากไม่มีปรับปรุงค่าถ่วงน้ำหนักตามที่ได้กล่าวมาในข้างต้น หรือกล่าวสรุปได้ว่าสถาปัตยกรรม RNN เกิดปัญหากับข้อมูลที่มีขนาดของลำดับ (Sequential) ยาวเกินไปนั่นเอง

2.18 Long Short-Term Memory (LSTM)

จากปัญหาของสถาปัตยกรรม RNN ที่ได้กล่าวไปข้างต้น มีการเสนอการใช้สถาปัตยกรรม LSTM โดย Sepp Hochreiter และ Juergen Schmidhuber



รูปที่ 2.6 สถาปัตยกรรม LSTM

(แหล่งที่มา: <https://medium.com/@sinart.t/long-short-term-memory-lstm-e6cb23b494c6>)

จากที่ได้กล่าวไปในข้างต้นแล้วว่า RNN เป็นเหมือนโครงข่ายประสาทเทียมที่มีหน่วยความจำ เพื่อบันทึกสถานะช่อก่อนหน้า และสำหรับ LSTM จะหน่วยความจำอยู่ภายในเช่นกัน แต่หลักสำคัญของสถาปัตยกรรม LSTM คือมีการทำงานที่บ่งบงได้ว่าเมื่อใดควรจะจดจำข้อมูลในหน่วยความจำ, ลืม (ลบ) ข้อมูลในหน่วยความจำหรืออนุญาตให้แบบจำลองรับข้อมูลอินพุตได้

2.18.1 การทำงานของ LSTM มีตัวแปรที่สำคัญดังนี้

1. เซลล์สถานะ (Cell State) คือตัวเก็บสถานะของเซลล์หน่วยความจำ (Memory Cell) ใน LSTM
2. Gate คือตัวที่ควบคุมการไหลของข้อมูล ซึ่งเปรียบเหมือนประตูที่คอยควบคุมว่าเมื่อไหร่ควรเปิดให้ข้อมูลไหลเข้า ไหลออก หรือไหลหายไปเลย (Forget)

ส่วนการทำงานต่าง ๆ ของ LSTM สามารถแบ่งได้ดังนี้

2.18.1.1 การลืม (Forget)

การลืม (Forget) คือการลบเซลล์สถานะเดิมออกไปเพื่อเตรียมพื้นที่รับข้อมูลใหม่ ซึ่งส่วนการทำงานนี้เป็นหน้าที่ตัดสินใจว่าจะลบหรือไม่ เป็นหน้าที่ของ Forget Gate

การสร้าง Forget Gate จะพิจารณาจากข้อมูลอินพุตประกอบกับสถานะช่อก่อนหน้าประกอบการตัดสินใจ โดยจะใช้ Sigmoid Function ในการตัดสินใจ

$$f_t = \sigma(W_{xf}x_t + W_{hf}h_{t-1} + b_f) \quad (2.10)$$

2.18.1.2 การเขียน (Write)

เมื่อมีข้อมูลอินพุตใหม่เข้ามา จะมีการพิจารณาอยู่ 2 ส่วนการทำงาน คือ

1. จะมีการปรับปรุงค่าเซลล์สถานะด้วย ข้อมูลอินพุตใหม่หรือไม่
2. หากมีการปรับปรุงค่า จะปรับปรุงด้วยค่าใดจึงจะเหมาะสม

พิจารณาการปรับปรุงค่าเซลล์สถานะ โดยมีการควบคุมส่วนงานนี้โดย Input Gate ซึ่งใช้ Sigmoid Function และใช้ข้อมูลอินพุตประกอบกับสถานะซ่อนก่อนหน้าประกอบการตัดสินใจว่าจะอนุญาตให้ปรับปรุงค่าหรือไม่

$$i_t = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + b_i) \quad (2.11)$$

ค่าที่เหมาะสมเมื่อมีการปรับปรุงค่าเซลล์สถานะ ซึ่งจะ Input Modulation gate เป็นตัวจัดการ โดยสมการก็จะคล้ายคลึงกับ Input Gate แต่จะใช้เป็น Tanh Function ในการคำนวณ ซึ่งค่าที่ได้จากการคำนวณเปรียบเสมือนว่าเป็น Cell State Candidate

$$g_t = \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c) \quad (2.12)$$

2.18.1.3 การปรับปรุงเซลล์สถานะ (Cell state)

จากการคำนวณของ Forget Gate, Input Gate และ Input Modulation Gate แล้ว เมื่อนำไปปรับปรุงเซลล์สถานะได้ตามสมการนี้

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t \quad (2.13)$$

จากสมการทางด้านฝั่งขวาในส่วนแรกของสมการจะเห็นได้ว่า ถ้า Forget Gate จะกำหนดให้ลบเซลล์สถานะเดิมทิ้ง เมื่อ f_t มีค่าเป็น 0 และมีการเก็บเซลล์สถานะ c_{t-1} ไว้ประกอบการพิจารณาการปรับปรุง เมื่อ f_t มีค่าเป็น 1

ในส่วนหลังของสมการ จะเป็นส่วนการทำงานปรับปรุงเซลล์สถานะจากข้อมูลใหม่ ซึ่งจะพิจารณาจากค่าที่ได้จาก Input Modulation Gate (g_t) ร่วมกับค่าผลลัพธ์จาก Input Gate (i_t) ว่ามีความเหมาะสมหรือไม่ หาก i_t มีค่าเป็น 1 จะใช้ค่า g_t ในการปรับปรุงเซลล์สถานะ แต่หาก i_t มีค่าเป็น 0 จะไม่ใช้ค่า g_t ที่มีมิติ (Dimension) นั้นในการปรับปรุงเซลล์สถานะ

2.18.1.4 การอ่าน (Read)

สำหรับผลลัพธ์ ณ เวลานั้น ๆ ของสถาปัตยกรรม LSTM จะมี 2 ส่วนคือ y_t และ h_t ในส่วนของ h_t จะมีการส่งต่อค่าไปยังเซลล์สถานะถัดไป ซึ่งจะมีการพิจารณาจะมีการส่งค่า h_t ไปยังเซลล์สถานะ h_{t+1} หรือไม่ หรือกล่าวว่ายอมเซลล์สถานะที่ $t=t+1$ อ่านค่าสถานะซ่อนหรือไม่ โดยในส่วนของการทำงานนี้ Output Gate ในการตัดสินใจ ซึ่งจะยังคงใช้สมการเดิมกับเหมือนกับการคำนวณ Forget Gate และ Input Gate ดังนี้

$$o_t = \sigma(W_{x^o}x_t + W_{h^o}h_{t-1} + b_o) \quad (2.14)$$

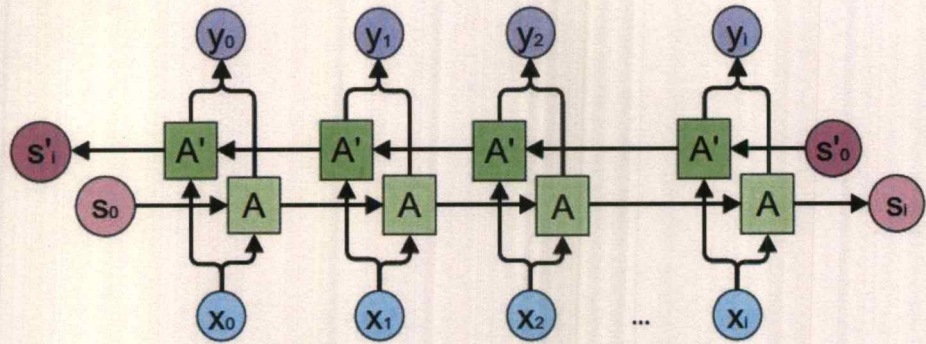
สมการที่ใช้ในการพิจารณาการส่งออกค่า h_t ไปยังเซลล์สถานะถัดไป คือ

$$h_t = o_t \odot \tanh(c_t) \quad (2.15)$$

จากสมการจะเห็นได้ว่าหาก Output Gate มีผลลัพธ์ค่า o_t เป็น 0 จะทำให้ได้ค่าผลลัพธ์ของ h_t มีค่าเป็น 0 ซึ่งมีความหมายว่าจะไม่ส่งค่าใด ๆ ออกไปไปยังเซลล์สถานะถัดไป และในขณะเดียวกันหาก o_t มีค่าเป็น 1 มีความหมายว่าจะคำนวณค่า h_t และส่งค่าผลลัพธ์ไปยังเซลล์สถานะถัดไป

2.19 Bi-directional LSTM [15]

Bi-directional LSTM หรือ LSTM แบบสองทิศทาง ถูกพัฒนาต่อยอดจากสถาปัตยกรรม LSTM ซึ่งจะมีการการทำงานเหมือนกับ LSTM ทุกประการ แต่เมื่อทำการจนถึงเซลล์สถานะสุดท้ายแล้วจะทำการคำนวณย้อนกลับไปยังข้อมูลอินพุตลำดับแรกดังรูป



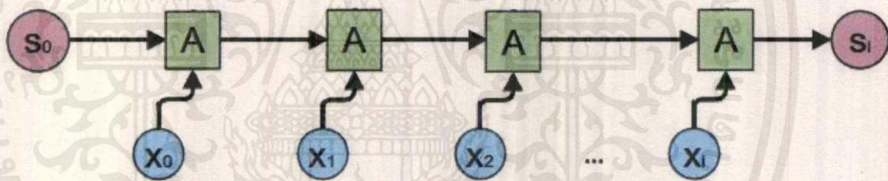
รูปที่ 2.7 สถาปัตยกรรม Bi-directional LSTM

(ที่มา: <https://towardsdatascience.com/understanding-bidirectional-rnn-in-pytorch-5bd25a5dd66>)

2.19.1 ขั้นตอนการทำงาน Bi-directional LSTM

สำหรับการขั้นตอนการดำเนินการของ Bi-directional LSTM มีดังนี้

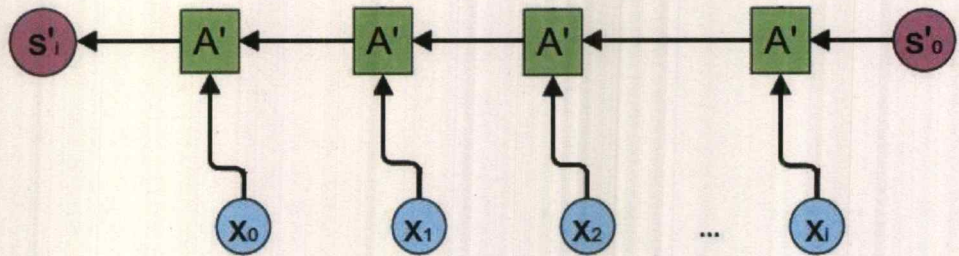
ขั้นตอนที่ 1: ดำเนินการเช่นเดียวกับสถาปัตยกรรม LSTM ทุกประการ



รูปที่ 2.8 Bi-directional LSTM ขั้นตอนไปข้างหน้า

(แก้ไขรูปจาก: <https://towardsdatascience.com/understanding-bidirectional-rnn-in-pytorch-5bd25a5dd66>)

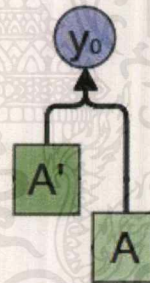
ขั้นตอนที่ 2: ดำเนินการเช่นเดียวกับสถาปัตยกรรม LSTM ทุกประการ แต่เริ่มต้นที่ข้อมูลอินพุตลำดับสุดท้ายไปยังข้อมูลอินพุตลำดับเริ่มต้น



รูปที่ 2.9 Bi-directional LSTM ขั้นตอนย้อนกลับ

(แก้ไขรูปจาก: <https://towardsdatascience.com/understanding-bidirectional-rnn-in-pytorch-5bd25a5dd66>)

ขั้นตอนที่ 3: นำค่าเวกเตอร์ของเซลล์สถานะที่มีข้อมูลลำดับเดียวกันมาต่อเรียงกัน (Concatenated) เช่น เวกเตอร์ $A = [1, 2, 3]$ และเวกเตอร์ $A' = [4, 5, 6]$ เมื่อดำเนินการ $A \text{ concat } A'$ จะได้ผลลัพธ์คือ $[1, 2, 3, 4, 5, 6]$

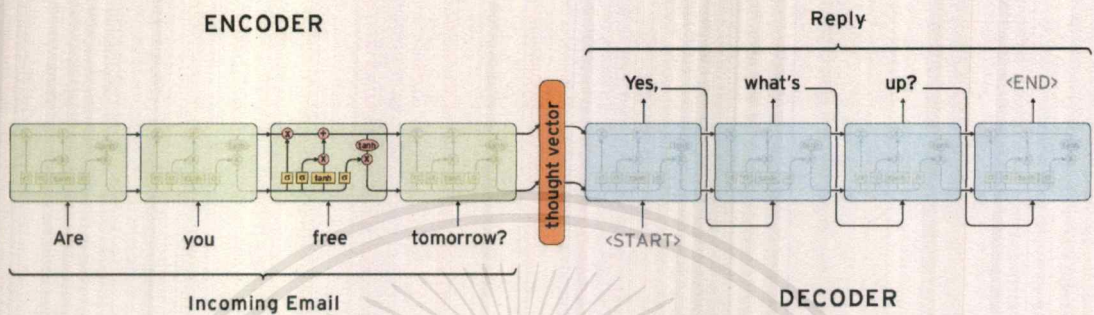


รูปที่ 2.10 การ Concatenate

(แก้ไขรูปจาก: <https://towardsdatascience.com/understanding-bidirectional-rnn-in-pytorch-5bd25a5dd66>)

2.20 แบบจำลอง Sequence-to-Sequence [16]

แบบจำลอง Sequence-to-Sequence เป็นการพัฒนาสถาปัตยกรรม LSTM เพื่อนำมาแก้ปัญหาที่มีข้อมูลอินพุตและเอาต์พุตในรูปแบบข้อมูลเป็นลำดับ (Sequential Data) เช่น การแปลภาษาด้วยเครื่อง, แชทบอท เป็นต้น



รูปที่ 2.11 แบบจำลอง Seq2Seq

(ที่มา: <https://towardsdatascience.com/sequence-to-sequence-tutorial-4fde3ee798d8>)

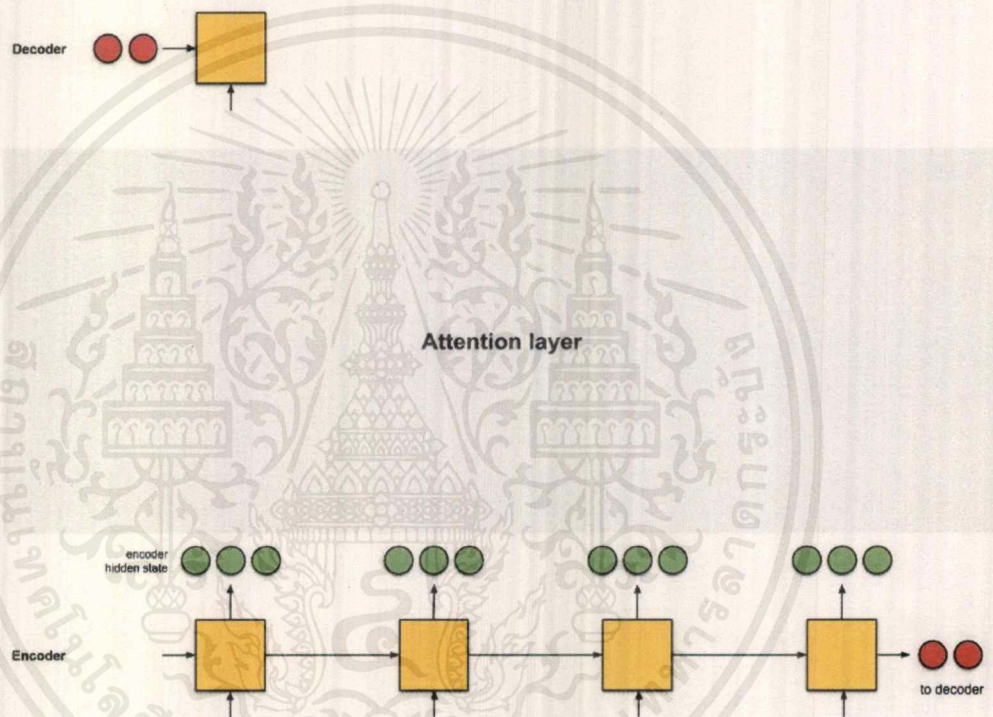
สำหรับการทำงานของแบบ Sequence-to-Sequence จะถูกแบ่งเป็น 2 ส่วนการทำงาน คือ Encoder และ Decoder ซึ่งมีรายละเอียดการทำงานดังนี้

1. Encoder จะทำงานร่วมกับข้อมูลอินพุต โดยมีการทำงานตามสถาปัตยกรรม LSTM ทุกประการ
2. Decoder จะทำงานร่วมกับข้อมูลเอาต์พุต จากรูป 2.11 จะเห็นได้ว่าในส่วนของเซลล์สถานะเริ่มต้นจะรับค่าเวกเตอร์มาจากส่วนการทำงาน Encoder และมีข้อมูลอินพุตเริ่มต้นคือคำว่า "<START>" จากนั้นจะดำเนินการทำนายผลลัพธ์และนำไปเป็นข้อมูลอินพุตในเซลล์ถัดไป โดยมีการดำเนินการต่าง ๆ ตามการทำงานของสถาปัตยกรรม LSTM ทุกประการ

2.21 Bahdanau Attention [17]

จากแบบจำลอง Sequence-to-Sequence ได้มีการพัฒนาเพื่อให้มีการทำงานที่ดีขึ้นและเหมาะสมในแต่ละปัญหา ซึ่งวิธีการที่ได้รับความนิยมคือทำการเพิ่ม Attention Layer เข้ามาในการประมวลผล อยู่ระหว่าง Encoder และ Decoder ดังรูป 2.12

Bahdanau Attention เป็น Attention Layer ที่ได้รับความนิยมอย่างกว้างขวางในการนำไปประยุกต์ใช้ในปัญหาต่าง ๆ เช่น การแปลภาษาด้วยเครื่อง, แชนบอท, การสรุปข้อความ เป็นต้น เนื่องจากมีความง่ายต่อการศึกษาและให้ผลลัพธ์ที่ดีในระดับหนึ่ง



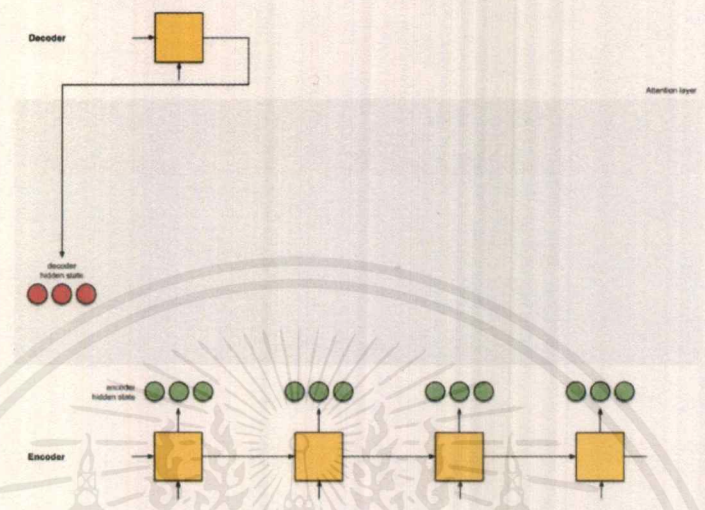
รูปที่ 2.12 Attention Layer

(ที่มา: <https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3>)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2.21.1 ขั้นตอนการทำงานของ Attention Layer

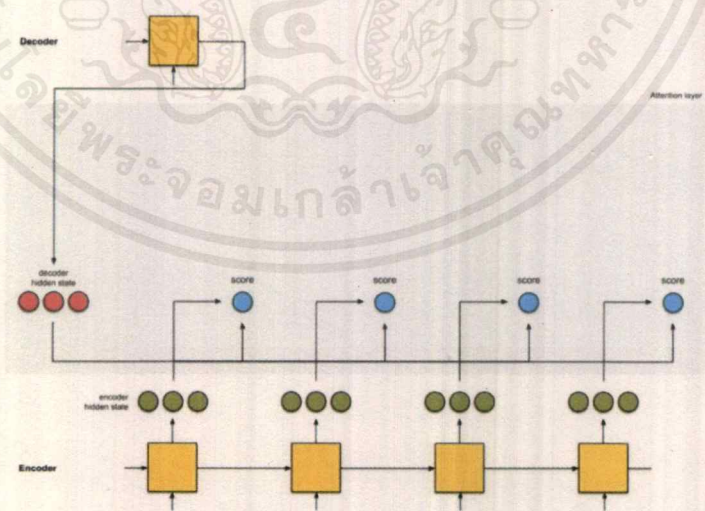
ขั้นตอนที่ 1: คำนวณค่าต่างๆจากเซลล์สถานะตามการทำการของสถาปัตยกรรม LSTM ทุกประการ



รูปที่ 2.13 Attention Layer การคำนวณขั้นตอนที่ 1

(ที่มา: <https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3>)

ขั้นตอนที่ 2: คำนวณคะแนนโดยใช้เซลล์สถานะของ Decode กับ เซลล์สถานะในแต่ละลำดับเวลา ดังรูป 2.14



รูปที่ 2.14 Attention Layer การคำนวณขั้นตอนที่ 2

(ที่มา: <https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3>)

โดยมีตัวอย่างการคำนวณดังนี้

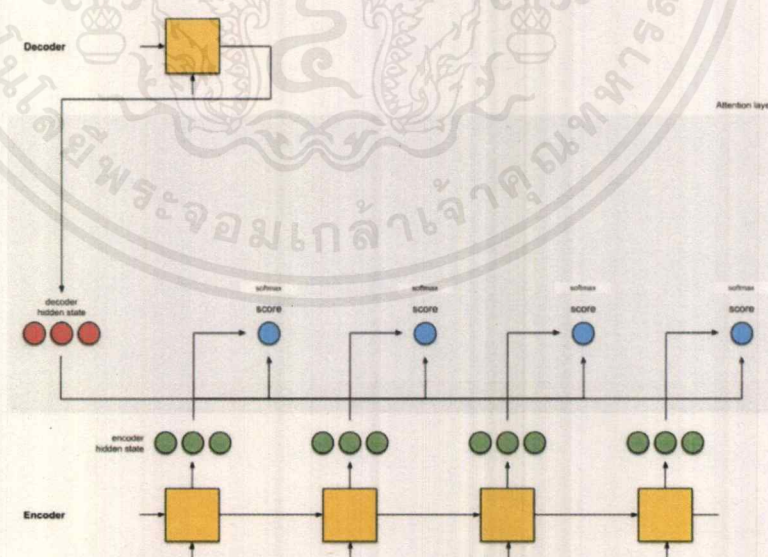
กำหนดให้ Decoder Hidden State = [10, 5, 10]

ดำเนินการคำนวณคะแนนของเซลล์สถานะในแต่ละลำดับด้วยการทำ Dot Product ระหว่าง Decoder Hidden State และ Encoder Hidden State จะได้ผลลัพธ์ดังนี้

ตารางที่ 2.8 แสดงผลลัพธ์การทำ Dot Product ระหว่าง Decoder Hidden State และ Encoder Hidden State

Encoder Hidden State	Score
[0, 1, 1]	$15 = (10 \times 0) + (5 \times 1) + (10 \times 1)$
[5, 0, 1]	$60 = (10 \times 5) + (5 \times 0) + (10 \times 1)$
[1, 1, 0]	$15 = (10 \times 1) + (5 \times 1) + (10 \times 0)$
[0, 5, 1]	$35 = (10 \times 0) + (5 \times 5) + (10 \times 1)$

ขั้นตอนที่ 3: นำผลลัพธ์จากการคำนวณในขั้นตอนที่ 2 ผ่านฟังก์ชัน Softmax โดยฟังก์ชันนี้จะทำการเลือกเวกเตอร์ที่มีคะแนนสูงสุด คือ 60 ดังนั้นในการดำเนินการขั้นตอนนี้จะได้ผลลัพธ์ คือ เวกเตอร์ [5, 0, 1]



รูปที่ 2.15 Attention Layer การคำนวณขั้นตอนที่ 3

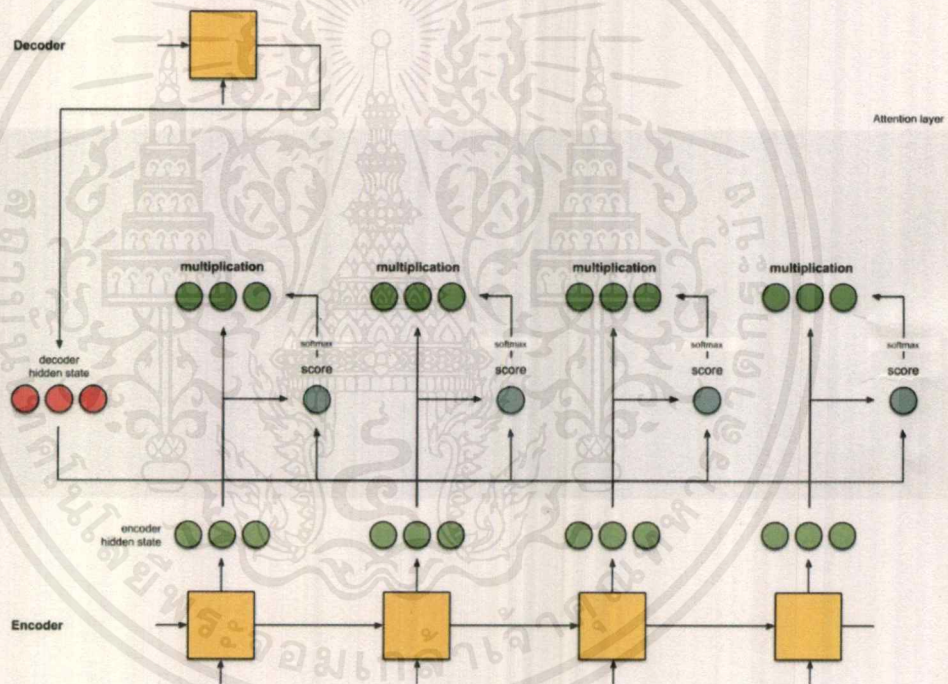
(ที่มา: <https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3>)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 2.9 ผลลัพธ์จากการผ่านฟังก์ชัน Softmax

Encoder Hidden State	Score	Score'
[0, 1, 1]	15	0
[5, 0, 1]	60	1
[1, 1, 0]	15	0
[0, 5, 1]	35	0

ขั้นตอนที่ 4: คูณคะแนนจากขั้นตอนที่ 4 กับค่าเวกเตอร์เซลล์สถานะในแต่ละลำดับเวลา



รูปที่ 2.16 Attention Layer การคำนวณขั้นตอนที่ 4

(ที่มา: <https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3>)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 2.10 ผลลัพธ์จากการคูณ Encoder Hidden State และ Score'

Encoder Hidden State	Score	Score'	Alignment
[0, 1, 1]	15	0	[0, 0, 0]
[5, 0, 1]	60	1	[5, 0, 1]
[1, 1, 0]	15	0	[0, 0, 0]
[0, 5, 1]	35	0	[0, 0, 0]

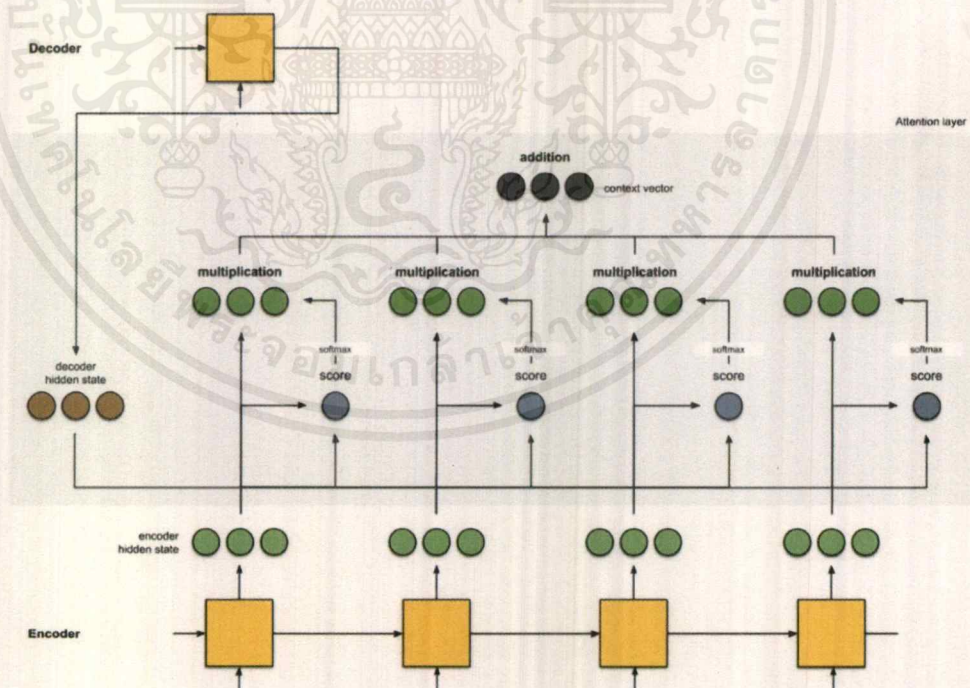
ขั้นตอนที่ 5: รวมผลลัพธ์เวกเตอร์ Alignment จากขั้นตอนที่สี่เข้าด้วยกัน โดยมีการรวมค่าตามตำแหน่งของเวกเตอร์ดังนี้

$$\text{ตำแหน่งที่ 1: } 0 + 5 + 0 + 0 = 5$$

$$\text{ตำแหน่งที่ 2: } 0 + 0 + 0 + 0 = 0$$

$$\text{ตำแหน่งที่ 3: } 0 + 1 + 0 + 0 = 1$$

ดังนั้นในขั้นตอนนี้จะได้ผลลัพธ์คือ [5, 0, 1] ซึ่งเรียกว่า เวกเตอร์บริบท (Context Vector)

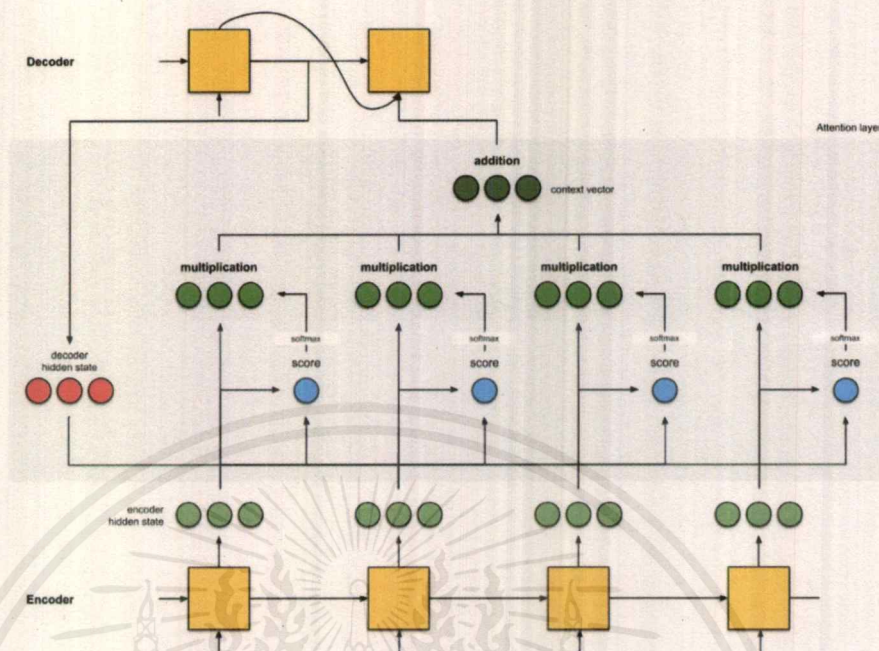


รูปที่ 2.17 Attention Layer การคำนวณขั้นตอนที่ 5

(ที่มา: <https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3>)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ขั้นตอนที่ 6: นำเวกเตอร์บริบทไปยังอินพุตในส่วนของ Decoder



รูปที่ 2.18 Attention Layer การคำนวณขั้นตอนที่ 6

(ที่มา: <https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3>)

2.22 เอกสารและงานวิจัยที่เกี่ยวข้อง

2.22.1 เอกสารที่เกี่ยวข้อง

2.22.1.1 Neural Abstractive Text Summarization with Sequence-to-Sequence Models [18]

รายงานฉบับนี้ได้รวบรวมแบบจำลองและเทคนิคต่าง ๆ ที่เกี่ยวข้องกับการสรุปข้อความประเภท Abstractive หรือที่เรียกว่า Neural Abstractive Text Summarizer (NATS) ตั้งแต่ปี พ.ศ. 2558 ถึงปี พ.ศ. 2561 ซึ่งแสดงให้เห็นถึงการให้ความสนใจจำนวนมากต่อการพัฒนาแบบจำลองการเรียนรู้ของเครื่องและเทคนิคต่าง ๆ เพื่อการสรุปข้อความให้มีความถูกต้องมากยิ่งขึ้น ทำให้ง่ายต่อการศึกษาและง่ายต่อการค้นหางานวิจัยที่เกี่ยวข้องกับปัญหาการสรุปข้อความ โดยมีการแสดงรายชื่อแบบจำลองการเรียนรู้ด้วยเครื่องหรือหรือสถาปัตยกรรมหรือเทคนิค ชื่อผู้วิจัย ซึ่งจัดเรียงลำดับตามช่วงเวลาในการพัฒนา

2.22.1.2 Abstractive Text Summarization [19]

เอกสารฉบับนี้กล่าวถึงที่มาและความสำคัญ ความรู้เบื้องต้น สถาปัตยกรรม และเทคนิคในการสรุปข้อความ โดยมีเนื้อหาโดยสรุปดังนี้

การสรุปข้อความถูกแบ่งออกเป็น 2 ประเภท คือ Extractive และ Abstractive

1. Extractive เปรียบเสมือนการอ่านหนังสือแล้วสกัดเนื้อหาสำคัญออกมาหรือคือการไฮไลต์ (Highlight)
2. Abstractive เปรียบเสมือนการอ่านและทำความเข้าใจเนื้อหาของข้อความและเขียนข้อความขึ้นมาใหม่ให้มีแต่เนื้อหาที่สำคัญ

เทคนิคที่มีความน่าสนใจและได้รับความนิยม มี 2 เทคนิค คือ 1. การนำแบบจำลอง Attention Mechanism มาร่วมประมวลผลกับแบบจำลอง Sequence-to-Sequence และ 2. การสรุปข้อความโดยการการเรียนรู้ด้วยเครื่องประเภทการเรียนรู้แบบเสริมแรง (Reinforcement Learning)

1. แบบจำลอง Attention Mechanism นิยมถูกใช้ในการแก้ปัญหาการแปลภาษาด้วยเครื่อง และมีการนำมาปรับใช้กับปัญหาอื่น ๆ ที่มีรูปแบบปัญหาที่ใกล้เคียง คือข้อมูลอินพุตและเอาต์พุตมีรูปแบบเป็นข้อมูลลำดับ ซึ่งในเอกสารฉบับนี้ได้กล่าวแนะนำแบบจำลอง Attention Mechanism ด้วย Bahdanau ซึ่งเป็นแบบจำลองที่สามารถศึกษาและทำความเข้าใจได้ง่ายในการเริ่มต้นเรียนรู้
2. การใช้การเรียนรู้แบบเสริมแรง (Reinforcement Learning) โดยกล่าวถึงการเรียนรู้แบบเสริมแรงโดยใช้ Policy Gradient-based โดยการประมวลผลในระดับคำ (Word Level) มีการคำนวณคะแนน (Reward) ด้วยตัวชี้วัด ROUGE สำหรับข้อสังเกตของวิธีการนี้คือ เป้าหมายของแบบจำลองไม่ควรถูกจำกัดให้แสดงผลลัพธ์ที่อ้างอิงความถูกต้องจากเพียงข้อความเป้าหมายเท่านั้น เนื่องจากในความเป็นจริงแล้ว ประโยคที่แตกต่างกันอาจจะมี ความหมายที่สื่อถึงสิ่งเดียวกันได้ หรือกล่าวคือในการสื่อ

ความหมายหนึ่งๆ สามารถถูกเขียนได้ในหลายรูปแบบ ซึ่งจากการคำนวณคะแนนด้วยตัวชี้วัด ROUGE จะทำให้ได้ผลลัพธ์คะแนนที่น้อย ซึ่งมีความไม่สมเหตุสมผลนั่นเอง

นอกจากนี้ยังมีการกล่าวถึงตัวชี้วัดในการวัดประสิทธิภาพในการสรุปข้อความของแบบจำลองการเรียนรู้ด้วยเครื่อง ซึ่งตัวชี้วัดที่ได้รับความนิยมสูงคือ ROUGE-L และ ROUGE-N

2.22.2 งานวิจัยที่เกี่ยวข้อง

2.22.2.1 Neural Machine Translation by Jointly Learning to Align and Translate (Bahdanau) [20]

จากปัญหาของแบบจำลอง Sequence-to-Sequence ที่จะนำเพียงเวกเตอร์บริบทหรือเวกเตอร์ของชั้นซ่อนลำดับสุดท้ายในส่วนของ Encoder ไปทำนายผลลัพธ์ในส่วนของ Decoder ซึ่งจากการดำเนินการนี้ทำให้เวกเตอร์ในชั้นซ่อนอื่น ๆ มีความสำคัญลดน้อยลงทำให้สูญเสียความถูกต้องของผลลัพธ์ในชั้นซ่อน Decoder ลำดับยาว (ลำดับที่อยู่ไกล)

จากปัญหาในข้างต้นในงานวิจัยนี้นำเสนอสถาปัตยกรรมสำหรับการแปลภาษาด้วยเครื่อง โดยมีการนำเวกเตอร์ของชั้นซ่อน Encoder ในทุก ๆ ลำดับไปประมวลผลร่วมกับส่วนการทำงาน Decoder

ผลลัพธ์จากงานวิจัยนี้ ทำให้ได้แบบจำลองที่สามารถสร้างประโยคที่มีความหมายคล้ายคลึงกัน กล่าวคือในการสื่อความหมายหนึ่งๆ สามารถถูกเขียนได้ในหลายรูปแบบ หรือมีการเขียนรูปแบบคำที่แตกต่างกันแต่สามารถสื่อความหมายเดียวกันได้

บทที่ 3

วิธีการดำเนินงาน

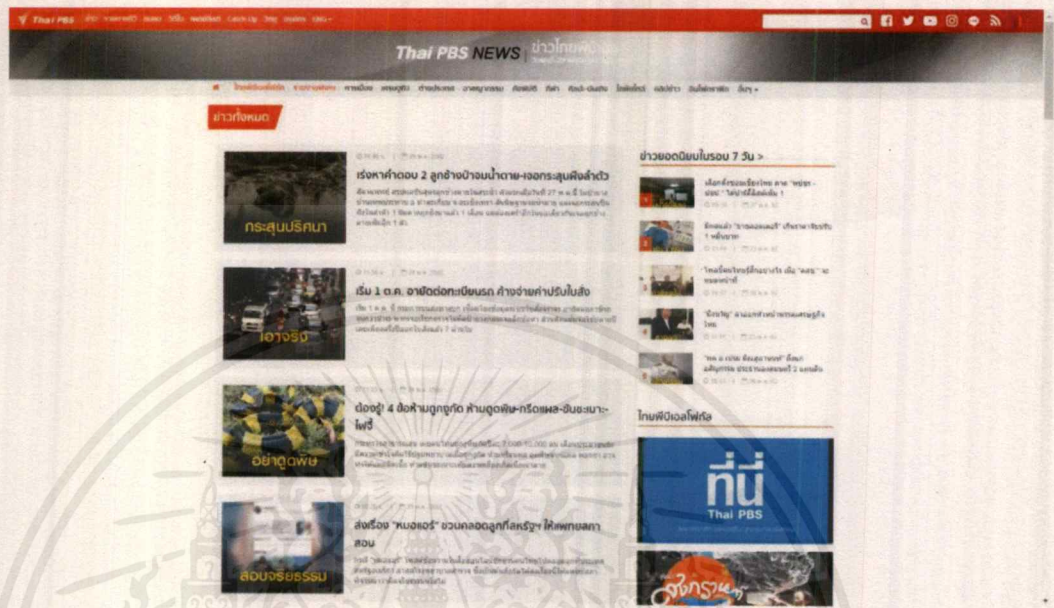
3.1 กระบวนการจัดเตรียมข้อมูล

ในการจัดเตรียมข้อมูลเพื่อใช้ในปัญหาพิเศษและเหมาะสมต่อการฝึกสอนแบบจำลองของเครื่อง มีขั้นตอนดังต่อไปนี้

1. เก็บรวบรวมข้อมูล
2. แปลงคำศัพท์ย่อ (Contractions)
3. ตัดคำ (Word Segmentation)
4. กำจัดคำหยุด (Stop Words)
5. กำจัดเครื่องหมายสัญลักษณ์
6. แปลงคำที่ไม่มีในพจนานุกรมเวกเตอร์เป็น “<UNK>”
7. เพิ่ม “<PAD>”, “<EOS>”, “<GO>” ในชุดข้อมูลเพื่อให้เหมาะสมต่อการสร้างแบบจำลองการเรียนรู้ของเครื่องประเภท Sequence-to-Sequence
8. แปลงคำศัพท์เป็นดัชนีของคำศัพท์ (Index)
9. แปลงดัชนีของคำศัพท์เป็นค่าเวกเตอร์ของคำศัพท์

3.1.1 การเก็บรวบรวมข้อมูล (Data Integration)

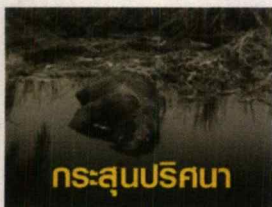
ข้อมูลที่นำมาใช้ในปัญหาพิเศษนี้รวบรวมมาจากเว็บไซต์สถานีโทรทัศน์ไทยพีบีเอส ในส่วนของเมนูข่าวทั้งหมด ซึ่งจะแสดงข่าวทั้งหมดทุกประเภท เรียงลำดับตามเวลาในการเขียนข่าว



รูปที่ 3.1 เว็บไซต์ ThaiPBS

(ที่มา: <https://news.thaipbs.or.th/archive>)


ข้อมูลที่นำมาสร้างแบบจำลองคือข้อมูลในส่วนของหัวข้อข่าวและข้อมูลข่าวอย่างสั้น จำนวน 29,170 ข่าว



16:46 น. | 29 พ.ค. 2562

เร่ร่อนคำตอบ 2 ลูกข้างป่าจมน้ำตาย-เจอกระสุนฝังลำตัว


สัตว์แพทย์ สรรพผลชันสูตรลูกข้างตายในสระน้ำ ตัวแรกเมื่อวันที่ 27 พ.ค. นี้ ในป่าข้างบ้านเทพประทาน อ.ท่าตะเียบน จ.ฉะเชิงเทรา สันนิษฐานจมน้ำตาย แต่เจอกระสุนปืนฝังในลำตัว 1 นัดคาดถูกยิงมาแล้ว 1 เดือน แต่ต้องตรวจอีกในมือเดียวกันเจอลูกข้างตายเพิ่มอีก 1 ตัว



15:58 น. | 29 พ.ค. 2562

เริ่ม 1 ต.ค. आयัดต่อทะเบียบนรค ค้างจ่ายค่าปรับใบสั่ง

เริ่ม 1 ต.ค. นี้ กรมการขนส่งทางบก เชื่อมโยงข้อมูลระบบในสั่งจราจร ภายัดต่อภาษีรถจนกว่าชำระ หากเจอเรียกตรวจไม่ติดป้ายของกรมเจออีกขอหา ส่วนคดีแค้นจ่อใบปลิวใบฝอยเพียงครึ่งมือออกใบสั่งแล้ว 7 ล้านใบ



15:33 น. | 29 พ.ค. 2562

ต้องรู้! 4 ข้อห้ามถูกขูด ห้ามดูดพิษ-กรีดแผล-ขันชะเนาะ-ไฟจี

กระทรวงสาธารณสุข เผยคนไทยถูกงูพิษกัดมีละ 7,000-10,000 คน เดือนประชาชนยังมีความเข้าใจผิดวิธีปฐมพยาบาลเมื่อถูกงูกัด ห้ามกรีดแผล ดูดพิษจากแผล ทอหมา ฉายาทำให้แผลติดเชื้อ ห้ามขันชะเนาะเพิ่มความเสี่ยงเกิดเนือเน่าตาย

รูปที่ 3.2 ส่วนของข้อมูลข่าวที่นำมาประมวลผล

3.1.2 แปลงคำศัพท์ย่อ (Contractions)

ในภาษาอังกฤษมีการลดรูปของคำอยู่จำนวนมาก เช่น it's, you're ซึ่งมีความหมายเหมือน it is และ you are ตามลำดับ ดังนั้นในขั้นตอนนี้จะดำเนินการเปลี่ยนแปลงคำเขียนย่อให้อยู่ในรูปคำเขียนเต็ม

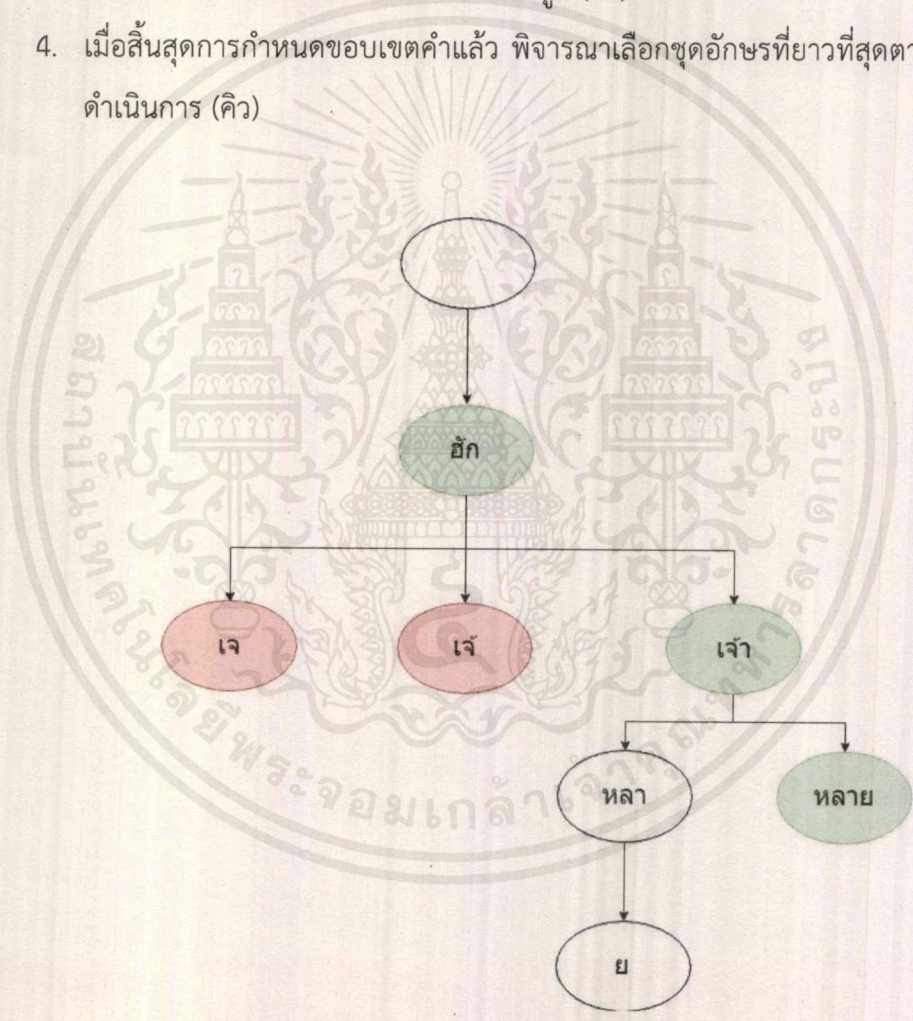
3.1.3 ตัดคำ (Word Segmentation)

ในปัญหาพิเศษนี้ใช้ตัวตัดคำ newmm [21] ซึ่งตัวตัดคำนี้ได้ถูกพัฒนาขึ้นโดยใช้เทคนิคการตัดคำแบบสอดคล้องมากที่สุดร่วมกับเทคนิคการตัดคำด้วยกลุ่มตัวอักษรไทย (Thai Character Cluster: TCC) ซึ่งมีหลักการทำงานดังนี้

1. ค้นหาความเป็นไปได้ทั้งหมดในการกำหนดขอบเขตคำ ด้วยการพิจารณาจากคลังคำศัพท์ว่ามีคำใดบ้างที่เป็นส่วนประกอบ (Subset) ของชุดอักษรที่ต้องการกำหนดขอบเขตคำ โดยเริ่มต้นที่อักษรตัวแรก

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2. พิจารณาตำแหน่งของผลลัพธ์ที่กำหนดขอบเขตคำจากขั้นตอนที่ 1 โดยคัดเลือกการกำหนดขอบเขตคำที่สอดคล้องกับกฎของ TCC
 - หากมีจำนวนผลลัพธ์จากการคัดเลือกมากกว่า 0 ดำเนินการจัดเก็บคำเหล่านั้นในโครงสร้างข้อมูลประเภทคิว (Queue)
 - หากจำนวนผลลัพธ์จากการคัดเลือกเท่ากับ 0 จะดำเนินการกำหนดขอบเขตคำตามกฎของ TCC และจัดเก็บชุดอักษรในโครงสร้างข้อมูลประเภทคิวเช่นเดียวกัน
3. ดำเนินการซ้ำในขั้นตอนที่ 1 โดยเริ่มต้นที่ตำแหน่งการค้นหาที่ตำแหน่งต่อจากผลลัพธ์ในขั้นตอนที่ 2 เรียงตามลำดับการจัดเก็บข้อมูล (คิว)
4. เมื่อสิ้นสุดการกำหนดขอบเขตคำแล้ว พิจารณาเลือกชุดอักขรที่ยาวที่สุดตามลำดับการดำเนินการ (คิว)



รูปที่ 3.3 ขั้นตอนการดำเนินการกำหนดขอบเขตคำด้วยเทคนิค newmm กับคำว่า “ฮักเจ้าหลาย”

จากรูปที่ 3.3 แสดงขั้นตอนการดำเนินการกำหนดขอบเขตคำด้วยเทคนิค newmm กับคำว่า “ฮักเจ้าหลาย” ดังนี้

1. พิจารณาความเป็นไปได้ในการกำหนดขอบเขตคำและได้ผลลัพธ์เพียง 1 คำ คือ “ฮัก”
2. พิจารณาความเป็นไปได้ในการกำหนดขอบเขตคำต่อจากคำว่า “ฮัก” และได้ผลลัพธ์ 3 คำ ได้แก่ “เจ”, “เจ้” และ “เจ้า”
3. พิจารณาผลลัพธ์ทั้ง 3 คำสอดคล้องตามกฎของ TCC หรือไม่ ซึ่งได้ผลลัพธ์ว่า “เจ”, “เจ้” ไม่สอดคล้องตามกฎ เพราะเหตุผลว่า หากกำหนดขอบเขตคำที่คำว่า “เจ” ในลำดับการดำเนินการต่อไป จะเริ่มต้นการค้นหาที่อักษร “ฮ” ซึ่งไม่ถูกต้องตามกฎของ TCC และหากกำหนดขอบเขตที่คำว่า “เจ้” จะเริ่มต้นการค้นหาที่อักษร “า” ไม่ถูกต้องตามกฎของ TCC ด้วยเช่นกัน ดังนั้นผลลัพธ์ในขั้นตอนนี้มีเพียง 1 ผลลัพธ์คือคำว่า “เจ้า”
4. พิจารณาความเป็นไปได้ในการกำหนดขอบเขตคำต่อจากคำว่า “เจ้า” และได้ผลลัพธ์ 2 คำ ได้แก่ “หลา” และ “หลาย”
5. พิจารณาความเป็นไปได้ในการกำหนดขอบเขตคำต่อจากคำว่า “หลาย” ซึ่งมีเพียงอักษรเดียวคือ “ย” ซึ่งไม่พบในคลังคำศัพท์ จึงดำเนินการกำหนดขอบเขตคำด้วย TCC และได้ผลลัพธ์คือ “ย”
6. ดำเนินการคัดเลือกผลลัพธ์ทั้งหมดโดยเริ่มพิจารณาจากโหนดราก (Root Node) และเลือกชุดอักษรที่มีความยาวมากที่สุดในแต่ละลำดับชั้น (Level) ของแผนภาพต้นไม้ จะได้ผลลัพธ์คือ “ฮัก|เจ้าหลาย”

คลังพจนานุกรมที่ใช้ในการในการตัดคำนำมาจากโมดูล PyThaiNLP [22] ร่วมกับคลังคำศัพท์จากโมดูล NLTK [23] รวมทั้งหมด 298,756 คำ และเพิ่มคำศัพท์ที่เหมาะสมกับปัญหาพิเศษนี้ ได้แก่ กยศ., ซีพี, ดิวตี้, ทีซีดีซีเพชบุค, บุโร, พุลมุน, มธ., รีเมค, ส.ส.ท., สตร., อีอีซี, อีโอดี, เอฟซี, แบงค็อก, แอชแท็ก, ไอซ์, “, ”

3.1.4 กำจัดคำหยุด (Stop Words)

ในขั้นตอนนี้จะทำการลบคำในส่วนของอินพุต (Input) เพื่อกำจัดคำที่ไม่มีความหมายออกไป เนื่องจากเป็นคำที่ไม่มีความสำคัญต่อการประมวลผล เช่น คำเชื่อม, ประติษฐานวิเศษณ์ เป็นต้น แต่ไม่กำจัดคำหยุดในเอาต์พุต (Output) เนื่องจากต้องการให้ผลลัพธ์ที่ได้มีความเป็นธรรมชาติ เช่น “กินข้าวกับปลา” หากกำจัดคำหยุดจะได้ผลลัพธ์คือ “กินข้าวปลา” โดยใช้คลังคำศัพท์ Stop words จาก NLTK และ PyThaiNLP

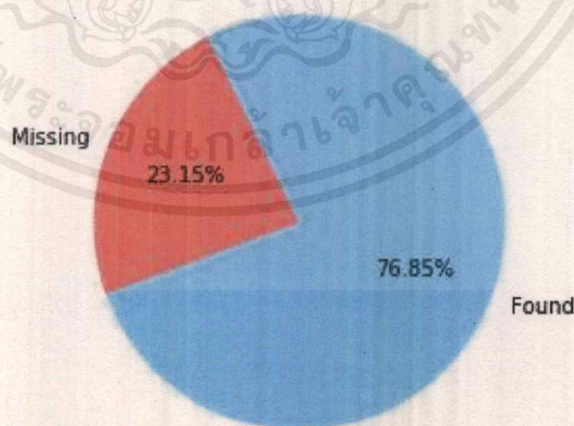
3.1.5 กำจัดเครื่องหมายสัญลักษณ์

เนื่องจากในบางเครื่องหมายไม่มีความสำคัญต่อการประมวลผลจึงทำการกำจัดออกไป เช่น {, }, ~, # เป็นต้น ในกรณีเครื่องหมาย “.” จะไม่ทำการกำจัดออกเนื่องมีการใช้เครื่องหมายดังกล่าว ในคำย่อต่าง ๆ เช่น i.e., U.S., กกต., กศน., ธ.ค. เป็นต้น

3.1.6 แปลงคำที่ไม่มีในพจนานุกรมเวกเตอร์เป็น “<UNK>”

สำหรับข้อมูลอินพุตของแบบจำลองจะเป็นค่าเป็นเวกเตอร์ของคำ โดยใช้พจนานุกรมเวกเตอร์ของ PyThaiNLP ซึ่งฝึกสอนด้วยแบบจำลอง Word2Vec ดังนั้นจึงต้องตรวจสอบก่อนว่ามีคำใดอยู่ในพจนานุกรมเวกเตอร์บ้าง และหากว่าอยู่ในรายการของพจนานุกรมเวกเตอร์จะทำการเปลี่ยนคำนั้น ๆ เป็น Unknow ซึ่งจะใช้สัญลักษณ์คือ “<UNK>” ในกรณีพบคำศัพท์คำหนึ่งๆ ในชุดข้อมูลจำนวนมากแต่ไม่พบในรายการของพจนานุกรมเวกเตอร์ จะทำการสุ่มค่าของเวกเตอร์ขึ้นมาเพื่อเพิ่มในพจนานุกรม ซึ่ง

Percentage of words found from vector dictionary



รูปที่ 3.4 แสดงอัตราส่วนของคำที่พบและไม่พบในพจนานุกรมเวกเตอร์

จากรูป 3.4 แสดงอัตราส่วนของคำที่พบและไม่พบในพจนานุกรมเวกเตอร์ ซึ่งมีจำนวน คำศัพท์ที่ไม่พบจำนวน 5,881 คำ คิดเป็นร้อยละ 23.15 และคำศัพท์ที่พบจำนวน 19,518 คำ คิด เป็นร้อยละ 76.85 จะได้ว่ามีคำในชุดข้อมูลทั้งหมดจำนวน 25,399 คำ

ในปัญหาพิเศษนี้กำหนดให้หากพบจำนวนคำหนึ่งๆปรากฏในชุดข้อมูลมากกว่า 20 ครั้งแต่ไม่ มีคำในพจนานุกรมเวกเตอร์ จะทำการสุ่มค่าของเวกเตอร์ให้คำนั้น ซึ่งมีจำนวน 359 คำ

3.1.7 เพิ่ม “<PAD>”, “<EOS>”, “<GO>” ในชุดข้อมูลเพื่อให้เหมาะสมต่อการสร้าง แบบจำลองประเภท Sequence-to-Sequence

2. “<EOS>”

เพิ่มสัญลักษณ์ “<EOS>” ในข้อความอินพุตและข้อความเอาต์พุตเพื่อให้เหมาะสม ต่อการฝึกสอนแบบจำลอง เช่น

[วิกฤต, ไฟป่า, ใน, แคลิฟอร์เนีย, เผา, ทำลาย, อาคาร, บ้าน, 1500, หลัง]
จะได้ [วิกฤต, ไฟป่า, ใน, แคลิฟอร์เนีย, เผา, ทำลาย, อาคาร, บ้าน,
1500, หลัง, <EOS>]

3. “<GO>”

เพิ่มสัญลักษณ์ “<GO>” ในข้อความเอาต์พุตเพื่อให้เหมาะสมต่อการฝึกสอน แบบจำลอง เช่น

[วิกฤต, ไฟป่า, ใน, แคลิฟอร์เนีย, เผา, ทำลาย, อาคาร, บ้าน, 1500, หลัง]
จะได้ [<GO>, ไฟป่า, ใน, แคลิฟอร์เนีย, เผา, ทำลาย, อาคาร, บ้าน, 1500, หลัง]

4. การ Padding

สำหรับข้อมูลอินพุตจำเป็นต้องทำให้ข้อมูลมีความยาวเท่ากันหรือมีจำนวนคำที่ เท่ากัน โดยกำหนดให้อิงตามจำนวนคำในข้อความข่าวที่มีความยาวมากที่สุด ซึ่งใน ปัญหาพิเศษนี้มีความยาวเท่ากับ 49 คำ ดังนั้นแล้วหากข้อความข่าวใดมีความยาวน้อย กว่า 49 จะทำการเพิ่มสัญลักษณ์ “<PAD>” เข้าไปในด้านหลัง เช่น

[วิกฤต, ไฟป่า, ใน, แคลิฟอร์เนีย, เผา, ทำลาย, อาคาร, บ้าน, 1500, หลัง,
<EOS>,<PAD>, <PAD>, <PAD>, <PAD>, ..., <PAD>]

จากตัวอย่างในข้างต้นข้อความข่าวมีจำนวน 10 คำ (ไม่นับรวม “<EOS>”) และทำ การเพิ่ม “<PAD>” อีก 39 คำ เพื่อให้ได้ความยาวของข้อความครบ 49 คำ

นอกจากนี้มีการดำเนินการกับข้อมูลเอาต์พุตเช่นเดียวกัน แต่ความยาวในการเพิ่ม “<PAD>” จะอิงตามจำนวนคำในข้อความสรุปข่าวที่มีความยาวมากที่สุด ซึ่งในปัญหาพิเศษนี้มีความยาวเท่ากับ 26 คำ

ดังนั้นในการดำเนินการขั้นตอนนี้จะได้ผลลัพธ์เช่นตัวอย่างดังต่อไปนี้

กำหนดให้ ข้อความอินพุตมีความยาวสูงสุด 20 คำ

 ข้อความเอาต์พุตมีความยาวสูงสุด 10 คำ

 ไม่มีรายการคำว่า “เป็ยน” ในพจนานุกรมเวกเตอร์

ข้อความก่อนดำเนินการ

ข้อความข่าว: [โปร, ช่าง, ธงชัย, ใจดี, โขว์, พอร์ม, เยี่ยม, คว่าแซมป์, กอล์ฟ, ยูโร, เป็ยน, ทัวร์, ฝรั่งเศส, ครอง, สำเร็จ]

ข้อความสรุปข่าว: [โปร, ช่าง, คว่าแซมป์, กอล์ฟ, ยูโร, เป็ยน, ทัวร์, ที่, ฝรั่งเศส]

ผลลัพธ์จากการดำเนินการในขั้นตอนนี้

ข้อความข่าว: [โปร, ช่าง, ธงชัย, ใจดี, โขว์, พอร์ม, เยี่ยม, คว่าแซมป์, กอล์ฟ, ยูโร, <UNK>, ทัวร์, ฝรั่งเศส, ครอง, สำเร็จ, <EOS>, <PAD>, <PAD>, <PAD>, <PAD>]

ข้อความสรุปข่าว: [โปร, ช่าง, คว่าแซมป์, กอล์ฟ, ยูโร, <UNK>, ทัวร์, ที่, ฝรั่งเศส, <EOS>, <PAD>]

3.1.8 แปลงคำศัพท์เป็นดัชนีของคำศัพท์ (Index)

ในปัญหาพิเศษนี้ใช้เครื่องมือสำหรับการเรียนรู้เชิงลึกคือ TensorFlow ซึ่งเป็นโมดูลในการฝึกสอนและพัฒนาแบบจำลองการเรียนรู้ของเครื่องในภาษาไพธอน (Python) และในส่วนงานนี้ TensorFlow จะช่วยในการจัดการทรัพยากรของหน่วยความจำเพื่อให้เพียงพอต่อการอินพุตข้อมูลดังเช่นตัวอย่างต่อไปนี้

[วิกฤต, ไฟป่า, ใน, แคลิฟอร์เนีย, เผา, ทำลาย, อาคาร, บ้าน, 1500, หลัง]

จากประโยคในข้างต้น หากต้องการนำข้อมูลไปฝึกสอนกับแบบจำลอง จะต้องดำเนินการแปลงข้อมูลในรูปแบบของข้อความให้เป็นข้อมูลในรูปแบบของเวกเตอร์ดังที่ได้กล่าวมาแล้ว ซึ่งคำหนึ่งๆนั้นจะถูกนำเสนอด้วยค่าเวกเตอร์จำนวน 300 มิติ ดังนั้นแล้วในประโยคข้างต้นจะมีมิติของข้อมูลคือ 10x300 ต่อข้อมูล 1 แถว หากมีข้อมูลจำนวน 1,000 แถว จะได้ว่า

1000x10x300 ซึ่งจากการดำเนินการนี้หากเป็นชุดข้อมูลที่มีปริมาณมากจะใช้พื้นที่หน่วยความจำมากและอาจไม่เพียงสำหรับผู้ที่มีทรัพยากรอย่างจำกัด

จากประโยคในข้างต้น กำหนดให้มีดัชนีเป็นดังนี้ [7, 9, 250, 88, 6, 5, 67, 65, 4500, 66] ซึ่งจะมีมิติของข้อมูลเพียง 1x10 ต่อข้อมูล 1 แถว หากมีข้อมูลจำนวน 1,000 แถว จะได้ว่ามีจำนวนมิติของข้อมูลเท่ากับ 1000x10

ด้วยเหตุนี้เอง TensorFlow จึงได้จัดเตรียมความช่วยเหลือในการแปลงคำศัพท์เป็นค่าเวกเตอร์ ซึ่ง TensorFlow จะรองรับข้อมูลประเภทดัชนีของคำศัพท์ ดังนั้นผู้พัฒนาโปรแกรมจึงจำเป็นต้องแปลงคำศัพท์เป็นดัชนีของคำศัพท์เพื่อให้ตรงตามที่ TensorFlow รองรับ

3.1.9 แปลงดัชนีของคำศัพท์เป็นค่าเวกเตอร์ของคำศัพท์

ในส่วนการทำงานนี้จะถูกนำเนินการในขั้นตอนการอินพุตข้อมูลไปยังแบบจำลอง ซึ่ง TensorFlow จะดำเนินการทั้งหมดให้โดยอัตโนมัติ สำหรับผู้พัฒนาโปรแกรม จะต้องจัดเตรียมให้ดัชนีของค่าเวกเตอร์ตรงตามดัชนีของคำศัพท์ตามที่ได้กล่าวไว้ในขั้นตอนที่ 8

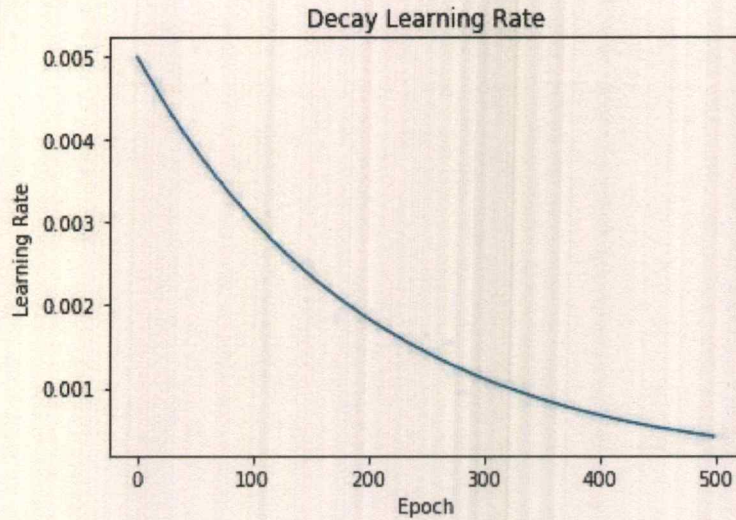
3.2 สร้างแบบจำลองการเรียนรู้ของเครื่อง

3.2.1 การกำหนดค่า Hyperparameter

- จำนวน Epoch เท่ากับ 500
- Optimizer คือ Adam [24] และ RMSProp [25] ซึ่งเป็น Optimizer ที่ได้รับความนิยมสูงเนื่องจากได้ปรับปรุงในส่วนด้อยของ Optimizer ตัวอื่น ๆ ให้ดียิ่งขึ้น
- Loss Function คือ Categorical Cross-Entropy
- ขนาดแบทช์ เท่ากับ 256
- จำนวนของ Units ใน LSTM เซลล์ เท่ากับ 256
- จำนวน LSTM เลเยอร์ เท่ากับ 2
- มีการคำนวณปรับค่าถ่วงน้ำหนักทุก ๆ 20 แบทช์
- Dropout ใน LSTM เซลล์ เท่ากับ 0 %, 20 % และ 35 %
- อัตราการเรียนรู้ (Learning Rate) กำหนดให้ค่าเริ่มต้นเท่ากับ 0.005 และมีการลดลงของอัตราการเรียนรู้ในแต่ละ Epoch คือ

$$\text{Learning Rate}_{\text{Epoch}(i)} = \text{Learning Rate}_{\text{Epoch}(i-1)} * 0.995 \quad (3.1)$$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 3.5 แสดงการลดลงของอัตราการเรียนรู้ที่คำนวณได้จากสมการ 3.1 เมื่อกำหนดให้อัตราการเรียนรู้เริ่มต้นเท่ากับ 0.005

3.2.2 การออกแบบจำลองการเรียนรู้ของเครื่อง

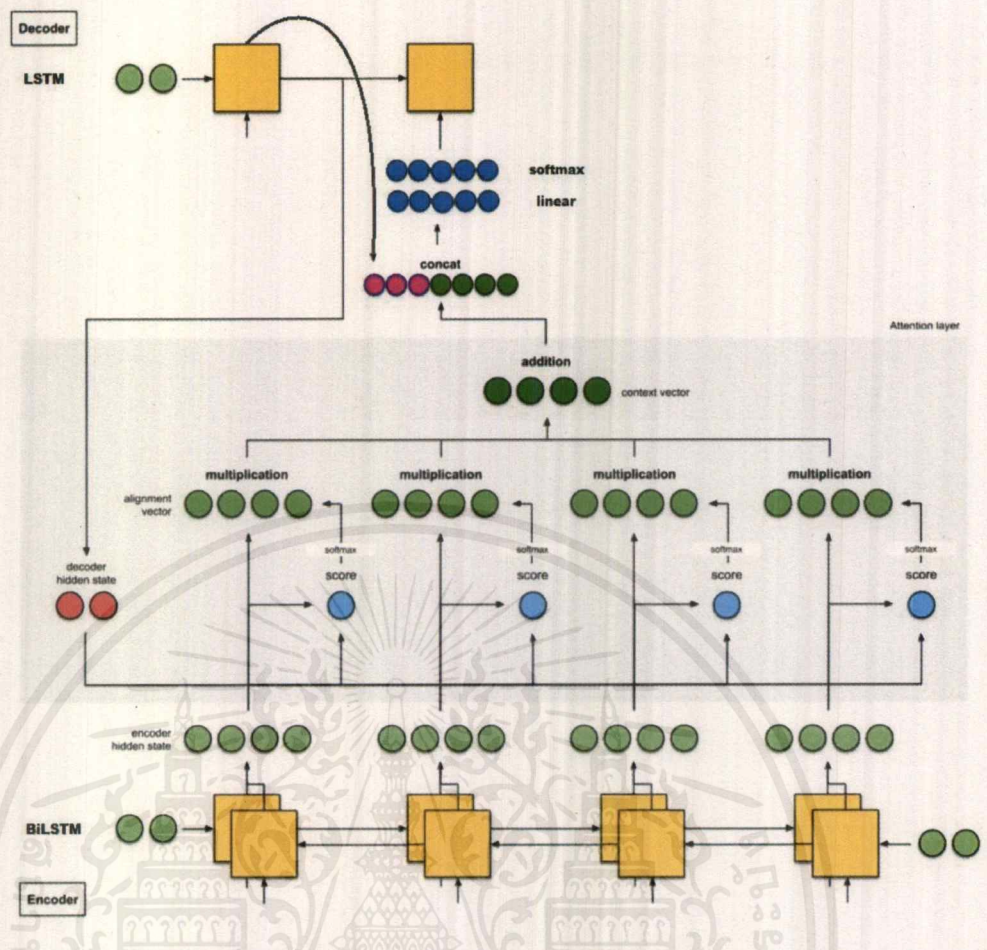
1. Encoder

- ใช้ Bi-directional LSTM
- มีข้อมูลอินพุต คือ เวกเตอร์ของคำ
- ใช้ Bahdanau ในการสร้าง Attention Layer

2. Decoder

- ใช้ LSTM
- ข้อมูลอินพุตคือเวกเตอร์ของคำ ซึ่งในการทำนายผลลัพธ์ว่าเป็นคำใดนั้นจะใช้โครงข่ายประสาทเทียม โดยมีรายละเอียดดังนี้

- ชั้นอินพุต (Input Layer) มีข้อมูลอินพุต คือเวกเตอร์ผลลัพธ์ของ Decoder ในเซลล์สถานะก่อนหน้าและเวกเตอร์ของผลลัพธ์มาเรียงต่อกัน (Concatenate) โดยจำนวนนิวรอนเท่ากับ 768
- ชั้นซ่อน (Hidden Layer) มีฟังก์ชันกระตุ้น คือฟังก์ชันเชิงเส้น (Linear Function) [26] โดยมีจำนวนนิวรอนเท่ากับ 19,878
- ชั้นเอาต์พุต (Output Layer) มีฟังก์ชันกระตุ้น คือฟังก์ชัน Softmax โดยมีจำนวนนิวรอนเท่ากับ 19,878



รูปที่ 3.6 แบบจำลอง Bahdanau

(แก้ไขรูปจาก: <https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3>)

3.3 การวัดประสิทธิภาพของแบบจำลองการเรียนรู้ของเครื่อง (Model Evaluation)

สำหรับปัญหาพิเศษนี้ใช้ตัวชี้วัด ROUGE-1, ROUGE-2 และ ROUGE-L ในการวัดประสิทธิภาพของแบบจำลองการเรียนรู้ของเครื่อง ซึ่งในแต่ละตัวชี้วัดประกอบไปด้วยค่าวัดประสิทธิภาพ ได้แก่ Recall, Precision และ F-measure

บทที่ 4

ผลการทดลองและอภิปรายผล

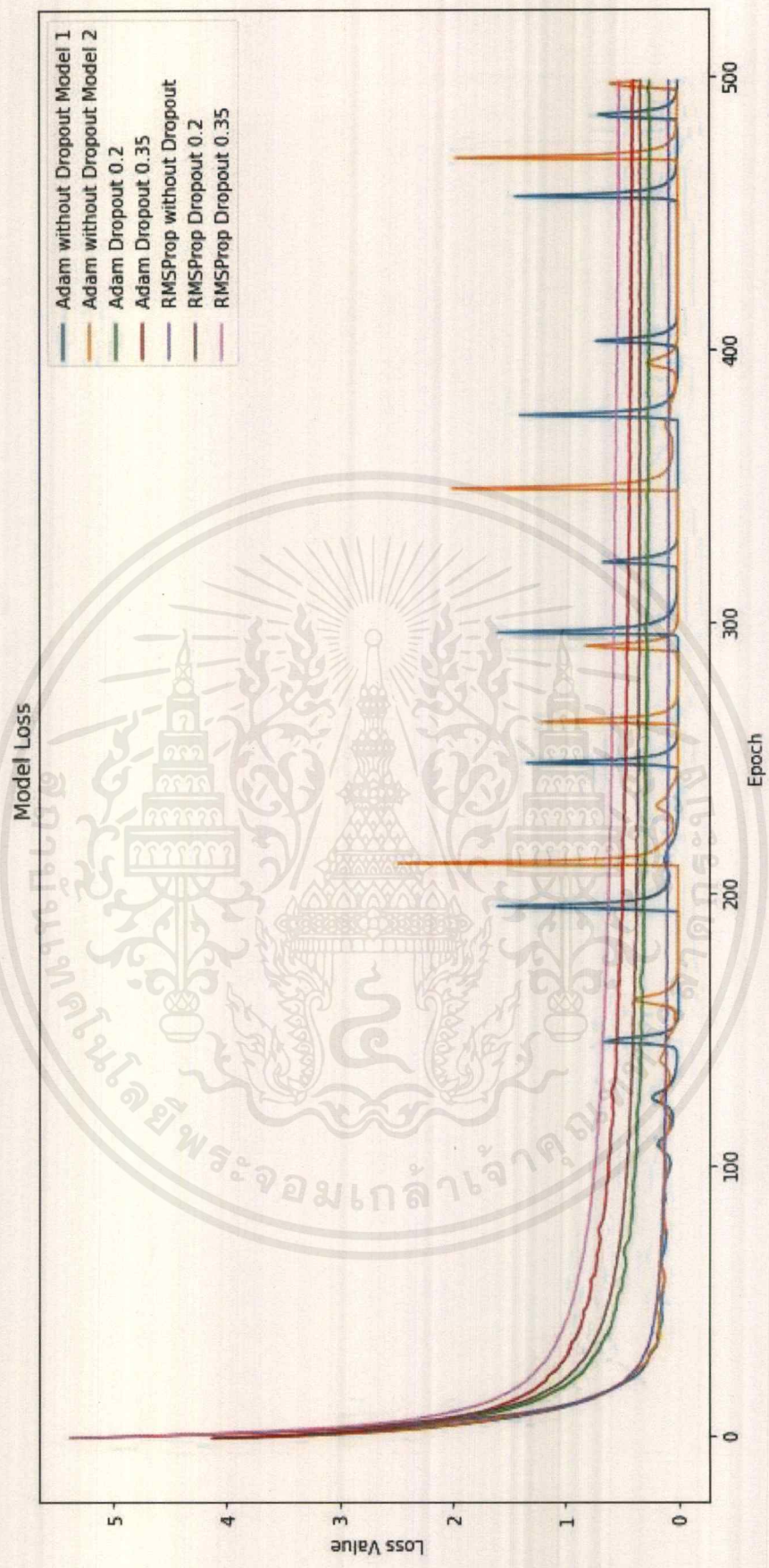
4.1 ผลการทดลอง

ในปัญหาพิเศษนี้มีการสร้างแบบจำลองทั้งหมด 7 แบบจำลอง ได้แก่

1. Adam without Dropout Model 1
2. Adam without Dropout Model 2
3. Adam Dropout 0.2
4. Adam Dropout 0.35
5. RMSProp without Dropout
6. RMSProp Dropout 0.2
7. RMSProp Dropout 0.35

ในการสร้างแบบจำลองโดยใช้ Adam Optimizer และไม่มีการ Dropout มีการสร้างแบบจำลอง 2 แบบจำลอง เนื่องจากพิจารณาพฤติกรรมของค่าความสูญเสีย (Loss) มีความผิดปกติ คือมีการเพิ่มขึ้นและลดลงในบางช่วง ได้แก่ แบบจำลองที่ 1: Adam without Dropout Model 1 และแบบจำลองที่ 2: Adam without Dropout Model 2

จากการทดลองสร้างแบบจำลองโดยการปรับเปลี่ยน Optimizer และค่าการ Dropout ได้ผลลัพธ์ค่าความสูญเสียดังรูปที่ 4.1



รูปที่ 4.1 แสดงค่าความสูญเสียของแบบจำลองในแต่ละ Epoch

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4.1 แสดงผลลัพธ์ค่าความถูกต้องโดยตัวชี้วัด ROUGE-1, ROUGE-2 และ ROUGE-L ประกอบด้วยค่าชี้วัด Recall, Precision และ F-measure

Evaluation metric	Model	Precision	Recall	F-measure
ROUGE-1	Adam without Dropout Model 1	0.2456	0.2171	0.2237
	Adam without Dropout Model 2	0.2518	0.2213	0.2285
	Adam Dropout 0.2	0.2104	0.2131	0.2047
	Adam Dropout 0.35	0.2171	0.2241	0.2138
	RMSProp without Dropout	0.2462	0.2075	0.2187
	RMSProp Dropout 0.2	0.2425	0.2247	0.2268
	RMSProp Dropout 0.35	0.2298	0.216	0.216
ROUGE-2	Adam without Dropout Model 1	0.0752	0.073	0.0715
	Adam without Dropout Model 2	0.0739	0.0714	0.07
	Adam Dropout 0.2	0.0564	0.0635	0.0574
	Adam Dropout 0.35	0.0605	0.0705	0.0627
	RMSProp without Dropout	0.0695	0.0645	0.0647
	RMSProp Dropout 0.2	0.0729	0.0749	0.0714
	RMSProp Dropout 0.35	0.0628	0.0666	0.0624
ROUGE-L	Adam without Dropout Model 1	0.2003	0.2031	0.1855
	Adam without Dropout Model 2	0.2027	0.2079	0.189
	Adam Dropout 0.2	0.1685	0.1989	0.165
	Adam Dropout 0.35	0.1714	0.2086	0.1698
	RMSProp without Dropout	0.1975	0.194	0.1809
	RMSProp Dropout 0.2	0.1964	0.2114	0.1878
	RMSProp Dropout 0.35	0.1825	0.203	0.1759

4.2 อภิปรายผลการทดลอง

จากตารางที่ 4.1 แสดงคะแนนถูกต้องในการสรุปเนื้อหาข่าว ซึ่งมีการคำนวณด้วยตัวชี้วัด ROUGE-1, ROUGE-2 และ ROUGE-L แบบจำลองที่ใช้ Adam Optimizer และไม่มีการ Dropout จะให้ผลลัพธ์ค่า Recall และ F-measure ดีกว่าแบบจำลองอื่น ๆ ในทุก ๆ ตัวชี้วัด และแบบจำลองที่ใช้ RMSProp Optimizer และมีการ Dropout 20 เปอร์เซ็นต์ จะมีค่า Precision ดีกว่าแบบจำลองอื่น ๆ ในทุก ๆ ตัวชี้วัด

จากรูปที่ 4.1 แสดงกราฟการค่าความสูญเสียในการเรียนรู้ในแต่ละ Epoch ของแบบจำลอง ซึ่งจากพฤติกรรมของค่าความสูญเสีย แสดงให้เห็นว่ามีค่าค่อนข้างคงที่ตั้งแต่ก่อนจะมีการเรียนรู้ครบจำนวน 100 Epoch เสียอีก เว้นแต่พฤติกรรมค่าความสูญเสียของแบบจำลอง ที่ใช้ Adam Optimizer และไม่มีการ Dropout จะมีการเพิ่มขึ้นของค่าความสูญเสียอย่างรวดเร็วและลดลงอย่างรวดเร็วในบางช่วง ซึ่งคาดการณ์ว่าอาจเป็นผลมาจากการมีข้อมูลที่เป็นข้อมูลรบกวน (Noise) หรืออาจมีการจัดเตรียมข้อมูลสำหรับการเรียนรู้ที่ไม่ถูกต้อง จึงได้มีการทดลองสร้างแบบจำลองจำนวน 2 แบบจำลองด้วยการกำหนดค่าต่าง ๆ ที่เหมือนกัน เพื่อศึกษาว่าจะเกิดพฤติกรรมค่าความสูญเสียเช่นนี้อีกหรือไม่ ซึ่งจากผลลัพธ์ว่าทั้ง 2 แบบจำลองมีพฤติกรรมค่าความสูญเสียในทำนองเดียว คือจะมีการเพิ่มขึ้นของค่าความสูญเสียอย่างรวดเร็วและลดลงอย่างรวดเร็วในบางช่วง และเมื่อนำแบบจำลองทั้ง 2 ไปวัดประสิทธิภาพการสรุปข้อความข่าว พบว่าแบบจำลองทั้ง 2 ให้ผลลัพธ์ค่าความถูกต้องที่ใกล้เคียงกัน

นอกจากนี้กราฟดังกล่าวยังแสดงให้เห็นว่าการ Dropout ที่มากขึ้น ทำให้เกิดค่าความสูญเสียที่มากขึ้น และเมื่อพิจารณาตารางที่ 4.2 จะเห็นได้ว่าเมื่อมีการ Dropout ที่มากขึ้นไม่ได้ทำให้มีค่าความถูกต้องในการสรุปเนื้อหาข่าวที่มากขึ้นตาม

บทที่ 5

สรุปผลการทดลองและข้อเสนอแนะ

5.1 สรุปผลการทดลอง

ปัญหาพิเศษนี้เป็นการศึกษาและนำเอาแบบจำลอง Bahdanau มาประยุกต์ใช้กับการสรุปข้อความข้ามภาษาไทย โดยข้อมูลที่นำมาใช้ในการสรุปข้อความข้ามคือข้อความข้ามจากสถานีโทรทัศน์ไทยพีบีเอส ซึ่งเป็นข้อมูลในส่วนของหัวข้อข่าวและข้อมูลข่าวอย่างสั้นจำนวน 29,170 ข่าว และใช้ตัวชี้วัด ROUGE ในการวัดประสิทธิภาพ

จากการทดลองสร้างแบบจำลองการเรียนรู้ด้วยเครื่องและวัดประสิทธิภาพได้ผลลัพธ์ว่ามีค่าความถูกต้องน้อย เนื่องจากมีจำนวนข้อมูลสำหรับการเรียนรู้ที่จำกัด ซึ่งไม่เหมาะสมต่อการเรียนรู้เชิงลึกที่ต้องการปริมาณข้อมูลจำนวนมาก แต่จากการพิจารณาผลลัพธ์การสร้างข้อความสรุปข่าวของแบบจำลองพบว่าผลลัพธ์มีความสนใจ คือแบบจำลองสามารถสร้างคำสรุปข่าวที่แตกต่างจากคำในเนื้อหาข่าวได้ โดยคำเหล่านั้นจะมีการสื่อความหมายที่เหมือนหรือคล้ายคลึงกัน

5.2 ข้อเสนอแนะ

1. หากต้องการใช้การเรียนรู้เชิงลึกในการสร้างแบบจำลองการเรียนรู้ของเครื่อง ควรจะมีข้อมูลจำนวนมากจึงจะเหมาะสม
2. ในการจัดสร้างแบบจำลองควรมีทรัพยากรในการประมวลผลที่สูง ซึ่งในปัจจุบันนิยมใช้หน่วยประมวลผลกราฟิกส์ เนื่องจากมีกำลังในการประมวลผลสูงและมีซอฟต์แวร์รองรับในการพัฒนาแบบจำลอง
3. ตัวตัดคำภาษาไทยยังไม่สามารถตัดคำศัพท์เฉพาะได้ดีเท่าที่ควร เช่น ชื่อของบุคคล ชื่อของสถานที่ เป็นต้น ซึ่งสามารถแก้ไขได้โดยการเพิ่มคำศัพท์ในพจนานุกรม ดังนั้นจึงควรมีการจัดสร้างพจนานุกรมดังกล่าวให้ครอบคลุมคำศัพท์เฉพาะมากที่สุด
4. ค่าเวกเตอร์ของคำในภาษาไทยยังจำเป็นต้องมีการพัฒนาต่ออีก เพื่อให้มีคลังคำศัพท์ที่มากขึ้นและเพื่อให้มีความถูกต้องของค่าเวกเตอร์ที่มากขึ้น
5. ในปัจจุบันมีการพัฒนาเทคนิคต่าง ๆ จำนวนมากสำหรับการแก้ปัญหาการสรุปข้อความ และยังมีเทคนิคในการแก้ปัญหาอื่น ๆ ที่มีความเกี่ยวข้องกัน เช่น การแปลภาษาด้วยเครื่อง ดังนั้นจำเป็นต้องมีการทดลองอีกมากสำหรับการนำมาประยุกต์ใช้กับภาษาไทย

เอกสารอ้างอิง

- [1] Thanaruk Theeramunkong, Virach Sornlertlamvanich, Thanasan Tanhermhong, Wirat Chinnan. 2000. "Character cluster based Thai information retrieval"
- [2] Tomas Mikolov. 2013. "Distributed Representations of Words and Phrases and their Compositionality"
- [3] "Categorical Cross-Entropy" [Online] Available from:
https://en.wikipedia.org/wiki/Cross_entropy Retrieved April 6, 2019
- [4] Chin-Yew Lin. 2004. "ROUGE: A Package for Automatic Evaluation of Summaries"
- [5] "Longest Common Subsequence Problem" [Online] Available from:
https://en.wikipedia.org/wiki/Longest_common_subsequence_problem
Retrieved June 1, 2019
- [6] "Deep learning" [Online] Available from:
https://en.wikipedia.org/wiki/Deep_learning Retrieved April 7, 2019
- [7] "Dropout (neural networks)" [Online] Available from:
[https://en.wikipedia.org/wiki/Dropout_\(neural_networks\)](https://en.wikipedia.org/wiki/Dropout_(neural_networks))
Retrieved April 7, 2019
- [8] "Long Short-Term Memory (LSTM)" [Online] Available from:
<https://medium.com/@sinart.t/long-short-term-memory-lstm-e6cb23b494c6>
Retrieved April 9, 2019
- [9] "Hyperbolic function" [Online] Available from:
https://en.wikipedia.org/wiki/Hyperbolic_function Retrieved April 12, 2019
- [10] "Rectifier (neural networks)" [Online] Available from:
[https://en.wikipedia.org/wiki/Rectifier_\(neural_networks\)](https://en.wikipedia.org/wiki/Rectifier_(neural_networks))
Retrieved April 12, 2019
- [11] "Sigmoid function" [Online] Available from:
https://en.wikipedia.org/wiki/Sigmoid_function Retrieved April 12, 2019

- [12] “Softmax function” [Online] Available from:
https://en.wikipedia.org/wiki/Softmax_function Retrieved April 12, 2019
- [13] “Backpropagation” [Online] Available from:
<https://en.wikipedia.org/wiki/Backpropagation> Retrieved April 13, 2019
- [14] “Vanishing gradient problem” [Online] Available from:
https://en.wikipedia.org/wiki/Vanishing_gradient_problem
 Retrieved April 13, 2019
- [15] “Understanding Bidirectional RNN in PyTorch” [Online] Available from:
<https://towardsdatascience.com/understanding-bidirectional-rnn-in-pytorch-5bd25a5dd66> Retrieved April 15, 2019
- [16] “Sequence-to-sequence Models” [Online] Available from:
<https://nlp.stanford.edu/~johnhew/public/14-seq2seq.pdf>
 Retrieved April 16, 2019
- [17] “Attn: Illustrated Attention” [Online] Available from:
<https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3>
 Retrieved April 21, 2019
- [18] Tian Shi, Yaser Keneshloo, Naren Ramakrishnan, Chandan K. Reddy. 2018.
 “Neural Abstractive Text Summarization with Sequence-to-Sequence Models”
- [19] “Abstractive Text Summarization” [Online] Available from:
<http://home.iitk.ac.in/~soumye/cs498a/report.pdf> Retrieved May 4, 2019
- [20] Dzmitry Bahdanau, Kyunghyun Cho, Yoshua Bengio. 2014. “Neural Machine Translation by Jointly Learning to Align and Translate”
- [21] “newmm algorithm” [Library] Available from:
<https://github.com/PyThaiNLP/pythainlp/blob/dev/pythainlp/tokenize/newmm.py>
- [22] “PyThaiNLP” [Library] Available from: <https://pythainlp.github.io/>
- [23] “NLTK” [Library] Available from: <https://www.nltk.org/>
- [24] Diederik P. Kingma, Jimmy Ba. 2017. “Adam: A Method for Stochastic Optimization”

- [25] “Neural Networks for Machine Learning – Lecture 6e rmsprop: Divide the gradient by a running average of its recent magnitude”
[Online] Available from:
https://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf
Retrieved May 15, 2019
- [26] “Linear predictor function” [Online] Available from:
https://en.wikipedia.org/wiki/Linear_predictor_function
Retrieved April 12, 2019





เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ภาคผนวก ก

รายการคำศัพท์ย่อ (Constructions)

ain't: am not	aren't: are not	can't: can not
can't've: cannot have	'cause: because	could've: could have
couldn't: could not	couldn't've: could not have	didn't: did not
doesn't: does not	don't: do not	hadn't: had not
hadn't've: had not have	hasn't: has not	haven't: have not
he'd: he would	he'd've: he would have	he'll: he will
he's: he is	how'd: how did	how'll: how will
how's: how is	i'd: i would	i'll: i will
i'm: i am	i've: i have	isn't: is not
it'd: it would	it'll: it will	it's: it is
let's: let us	ma'am: madam	mayn't: may not
might've: might have	mightn't: might not	must've: must have
mustn't: must not	needn't: need not	shan't: shall not
she'd: she would	she'll: she will	she's: she is
should've: should have	that's: that is	there'd: there had
there's: there is	they'd: they would	they'll: they will
they're: they are	they've: they have	wasn't: was not
we'd: we would	we'll: we will	we're: we are
we've: we have	weren't: were not	what'll: what will
what're: what are	what's: what is	what've: what have
where'd: where did	where's: where is	who'll: who will
who's: who is	won't: will not	wouldn't: would not
you'd: you would	you'll: you will	you're: you are
%: เปอร์เซ็นต์	\xa0: (ช่องว่าง 1 ช่อง)	
oughtn't: ought not	shan't: shall not	
shouldn't: should not	that'd: that would	

ภาคผนวก ข

ตัวอย่างผลลัพธ์การสรุปข้อความข่าว

ตัวอย่างการสรุปข้อความข่าวของแบบจำลองต่าง ๆ โดยมีลำดับผลลัพธ์ดังนี้

1. Adam without Dropout Model 1
2. Adam without Dropout Model 2
3. Adam Dropout 0.2
4. Adam Dropout 0.35
5. RMSProp without Dropout
6. RMSProp Dropout 0.2
7. RMSProp Dropout 0.35

สวนดุสิตโพลเผยประชาชนรู้สึกเบื่อหน่ายการจัดตั้งรัฐบาล โดยมองว่าเกิดจากผลประโยชน์ทางการเมืองไม่ลงตัว-แย่งชิงเก้าอี้รัฐมนตรี จึงอยากให้ยึดผลประโยชน์ของประเทศชาติเป็นหลัก ขณะที่นักโพลเผยสำรวจ 10 สาเหตุที่ทำให้พรรคประชาธิปัตย์ "แพ้เลือกตั้ง-แตกแยก"

1. โพลนิด้าช่วยร้องรัฐบาลแก้ปัญหาอดีต|นายทหารเกษียณ
2. โพลชี้ผู้ประกันตน|12ปี|เหตุแพ้เลือกตั้ง
3. โพลชี้ต้องมีผลกระทบ|คนวิจารณ์|เลือกตั้ง
4. โพลชี้อนาคตต้องเป็น|ของทักษิณ|หยุดทุจริต|จากเหตุเกษียณ
5. ผลโพลโพลหลัง|ผลทดสอบ|กลุ่มการเมือง
6. โพลชี้จับ|10|เก้าอี้|กลุ่มการเมือง
7. เลือกตั้ง|2562|เช็ก|ด่วน|ปชช.|ยังไม่ล้มเหลว|ทางการเมือง

เจ้าหน้าที่ท่าอากาศยานสุวรรณภูมิ ควบคุมตัวชาวจีน 2 คน พร้อมของกลาง เป็นทรัพย์สินที่
ขโมยจากผู้โดยสารต่างชาติระหว่างเดินทางเข้าไทย เจ้าหน้าที่เตรียมเสนอเพิกถอนการอนุญาตให้อยู่
ต่อในราชอาณาจักร พร้อมทั้งดำเนินคดี หลังเห็นว่ามีพฤติกรรมเป็นภัยสังคม

1. คุมตัวชาวจีนในสุวรรณภูมิเร่งเรือตัววันแรก
2. รวบผู้เข้าออฮอลลิวูดเพิ่มอำนาจผู้รุกในไทย
3. จนท.ที่ใช้เหยี่ยวโดยสารเหยื่อถูกเรียกแบะเจียะ
4. จับชาวจีนประสานชาวจีนต้องออกจากไทยจับเพิ่ม
5. นักดำน้ำเกาหลีเหนือเตรียมออกหมายจับเพิ่มคนก่อนขโมยทรัพย์สิน
มาแกะรอย
6. จนท.เร่งแก๊งหลังขโมยกระเป๋าในไทย
7. จนท.จับชาวจีนค้าคำสั่งลักทรัพย์ในสนามบินสุวรรณภูมิ

พล.อ.สุรยุทธ์ จุลานนท์ รักษาการประธานองคมนตรี เป็นผู้แทนพระองค์ ไปในการบำเพ็ญ
พระราชกุศลทักษิณานุปทาน 7 วัน พระราชทานศพ พล.อ.เปรม ติณสูลานนท์ อดีตประธาน
องคมนตรีและรัฐบุรุษ โดยมีบุคคลสำคัญและประชาชนเข้าร่วมพิธี

1. จากพลอ.เปรมชัยนั่งบัญชาให้ล.อ.เปรมเป็นประธาน
2. ล.อ.เปรมนั่งหัวหน้าพรรค300ตำแหน่ง300เปอร์เซ็นต์
3. พลอ.เปรมเปิดพิธีแสดงความอาลัยอดีตประธานองคมนตรี
4. ล.อ.เปรมเชื่อเป็นนักธุรกิจองคมนตรีทำบุญนำไปกลับไปพระราชกุศล
5. องคมนตรีวิงเปรมลาออกพลอ.เปรมเป็นพลอ.เปรมเป็นพลอ.เปรม
6. พลอ.เปรมชัยเผยประชาชนร่วมแสดงขอให้บรรจุศพล.อ.เปรม
7. ล.อ.เปรมรับพิธีเป็นประธานองคมนตรี

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

คปภ.แจ้งสิทธิคำรักษาพยาบาลและค่าสินไหมทดแทนตาม พ.ร.บ. ให้กับผู้บาดเจ็บและผู้เสียชีวิตแล้ว กรณีรถเบนซ์พุ่งชน บริเวณถนนงามวงศ์วาน พร้อมเร่งติดตามบริษัทประกันภัย

1. ติดตาม|ข้อมูล|ชรก. |คาด|ว่า|เร่ง|ติดตาม|ผู้|ก่อ|เหตุ|ใน|อดีต|ผู้|ใหญ่|บ้าน
2. คปภ.|เตรียม|ประเมิน|20|ม.ค.|นี้|เร่ง|ค้นหา|รถ|เสีย|ดัง|เบ|จ.|งาม
3. ตั้ง|ข้อ|กล่าวหา|สาเหตุ|รถ|เบนซ์|ออก|โฉนด|ที่ดิน|ถ.|พหลโยธิน
4. คปภ.|แจ้ง|ข้อ|กล่าวหา|คน|แอบ|อ้าง|ทำ|ญาติ|ออก|เบรค|การ|คุม|เข้ม|เหยี่ยว|ยา|3000|หมื่น|คน
5. คปภ.|จัด|กรม|ศึ|ไกร|เพชร|บูรณ|เสียชีวิต|ตร.|เร่ง|ตรวจสอบ|เพิ่ม|ตาม|ควบคุม|น้อง|ตาย|ควบคุม
6. คปภ.|เผย|เรียกร้อง|ตรวจ|เค|ดา|ปา|เข้า|ก่อสร้าง|อิตา|เลียน
7. คปภ.|เรียกร้อง|รอ|การใช้|สิทธิ์|นี้

เกิดเหตุผู้ก่อเหตุ 2 คนซึ่งจักรยานยนต์ ปาระเบิดไม่ทราบชนิดใส่ป้อมตำรวจที่ อ.บันนังสตา จ.ยะลา วัยรุ่นชาย อายุ 17 ปี โดนลูกหลงบาดเจ็บ 1 คน

1. ระทึก|พลัดตก|บาดเจ็บ|ก่อน|ปา|ระเบิด|ใส่|ประทัด
2. รถกระบะ|ใหม่|ขึ้น|จาก|ปา|แก๊ส|ใส่|นัก|รบ|ตร.|ส่ง|เว|บาดเจ็บ|คน
3. เกิดเหตุ|ขว้าง|ระเบิด|ใส่|บลู|ปา|แก้ว|ใส่|ฟิลิป|ปินส์|พุ่ง|เจ็บ|สาหัส
4. มือ|ปา|ระเบิด|ใส่|หัว|ตร.|ใส่|จาก|เหตุ|ใส่|ใส่|บาง|สะพาน|ตร.|เสียชีวิต|คน|ใส่|เร็ว
5. พิช|ปา|ถล่ม|ใส่|ระเบิด|ใน|จ.|ยะ|ไม่มี|เป็น|18|สาหัส
6. หลักฐาน|ชัด|อีก|ศพ|แก๊ง|ปา|ระเบิด|ใส่|รถกระบะ|ใส่|พลิก|คว่ำ|ปา|บึก|เจ็บ
7. เกิดเหตุ|ปา|ระเบิด|ใส่|ร้าน|สี|ปา|โฮ|ใส่|ร้าน|ตำรวจ|จ.|กาญจนบุรี|บาดเจ็บ|คน

ภาคผนวก ค

โค้ดที่ใช้ในการสร้างแบบจำลอง

1. โค้ดสร้างแบบจำลอง Adam without Dropout Model 1

```

1 from tensorflow.layers import Dense
2 from tensorflow.nn.rnn_cell import BasicLSTMCell
3 import numpy as np
4 import tensorflow as tf
5 import re
6 import time
7 import json
8
9 print('TensorFlow Version: {}'.format(tf.__version__))

```

TensorFlow Version: 1.13.1

```

1 import json
2
3
4 def read_json(fname, key_int=False):
5     with open(fname, 'r') as file:
6         data = file.read()
7         json_data = json.loads(data)
8
9         if not key_int:
10            return json_data
11
12            json_data = {int(key): value for key, value in json_data.items()}
13            return json_data
14
15
16 def to_json(fname, data):
17     with open(fname, 'w') as file:
18         json_data = json.dumps(data)
19         file.write(json_data)
20
21
22 def to_txt(fname, data):
23     with open(fname, 'w', encoding='utf-8') as file:
24         data = [str(d) for d in data]
25         file.write("\n".join(data))

```

Load value

```

1 vocab_to_int = read_json('../vocab_to_int.json')
2 int_to_vocab = read_json('../int_to_vocab.json', key_int=True)
3 word_embedding_matrix = np.load('../word_embedding_matrix.npy') # load
4
5 # dataset
6 dataset = read_json('../dataset.json')
7 X_train = dataset['X_train']
8 X_test = dataset['X_test']
9 y_train = dataset['y_train']
10 y_test = dataset['y_test']

```

```

1 print('Train: ', len(X_train))
2 print('Test: ', len(X_test))

```

Train: 21393

Test: 2377

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Building the Model

```

1 def model_inputs():
2     """Create palceholders for inputs to the model"""
3
4     input_data = tf.placeholder(tf.int32, [None, None], name='input')
5     targets = tf.placeholder(tf.int32, [None, None], name='targets')
6     lr = tf.placeholder(tf.float32, name='learning_rate')
7     keep_prob = tf.placeholder(tf.float32, name='keep_prob')
8     summary_length = tf.placeholder(tf.int32, (None,)), name='summary_length')
9     max_summary_length = tf.reduce_max(summary_length, name='max_dec_len')
10    text_length = tf.placeholder(tf.int32, (None,)), name='text_length')
11
12    return input_data, targets, lr, keep_prob, summary_length, max_summary_length, text_length

```

```

1 def process_encoding_input(target_data, vocab_to_int, batch_size):
2     """Remove the last word id from each batch and concat the <GO> to the beginning of each batch"""
3
4     ending = tf.strided_slice(target_data, [0, 0], [batch_size, -1], [1, 1])
5     dec_input = tf.concat([tf.fill([batch_size, 1], vocab_to_int['<GO>']), ending], 1)
6
7     return dec_input

```

```

1 def encoding_layer(rnn_size, sequence_length, num_layers, rnn_inputs, keep_prob):
2     """Create the encoding layer"""
3
4     for layer in range(num_layers):
5         with tf.variable_scope('encoder_{}'.format(layer)):
6             cell_fw = tf.contrib.rnn.LSTMCell(rnn_size,
7                 initializer=tf.random_uniform_initializer(-0.1, 0.1, seed=2))
8             cell_fw = tf.contrib.rnn.DropoutWrapper(cell_fw,
9                 input_keep_prob = keep_prob)
10
11            cell_bw = tf.contrib.rnn.LSTMCell(rnn_size,
12                initializer=tf.random_uniform_initializer(-0.1, 0.1, seed=2))
13            cell_bw = tf.contrib.rnn.DropoutWrapper(cell_bw,
14                input_keep_prob = keep_prob)
15
16            enc_output, enc_state = tf.nn.bidirectional_dynamic_rnn(cell_fw,
17                cell_bw,
18                rnn_inputs,
19                sequence_length,
20                dtype=tf.float32)
21
22            # Join outputs since we are using a bidirectional RNN
23            enc_output = tf.concat(enc_output, 2)
24
25            return enc_output, enc_state

```

```

1 def training_decoding_layer(dec_embed_input, summary_length, dec_cell, initial_state, output_layer,
2     vocab_size, max_summary_length):
3     """Create the training logits"""
4
5     training_helper = tf.contrib.seq2seq.TrainingHelper(inputs=dec_embed_input,
6         sequence_length=summary_length,
7         time_major=False)
8
9     training_decoder = tf.contrib.seq2seq.BasicDecoder(dec_cell,
10         training_helper,
11         initial_state,
12         output_layer)
13
14     training_logits, _, _ = tf.contrib.seq2seq.dynamic_decode(training_decoder,
15         output_time_major=False,
16         impute_finished=True,
17         maximum_iterations=max_summary_length)
18
19     return training_logits

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

1 def inference_decoding_layer(embeddings, start_token, end_token, dec_cell, initial_state, output_layer,
2                               max_summary_length, batch_size):
3     """Create the inference logits"""
4
5     start_tokens = tf.tile(tf.constant([start_token], dtype=tf.int32), [batch_size], name='start_tokens')
6
7     inference_helper = tf.contrib.seq2seq.GreedyEmbeddingHelper(embeddings,
8                                                                start_tokens,
9                                                                end_token)
10
11    inference_decoder = tf.contrib.seq2seq.BasicDecoder(dec_cell,
12                                                       inference_helper,
13                                                       initial_state,
14                                                       output_layer)
15
16    inference_logits, _, _ = tf.contrib.seq2seq.dynamic_decode(inference_decoder,
17                                                             output_time_major=False,
18                                                             impute_finished=True,
19                                                             maximum_iterations=max_summary_length)
20
21    return inference_logits

```

```

1 def decoding_layer(dec_embed_input, embeddings, enc_output, enc_state, vocab_size, text_length, summary_length,
2                   max_summary_length, rnn_size, vocab_to_int, keep_prob, batch_size, num_layers):
3     """Create the decoding cell and attention for the training and inference decoding layers"""
4
5     for layer in range(num_layers):
6         with tf.variable_scope('decoder_{}'.format(layer)):
7             lstm = tf.contrib.rnn.LSTMCell(rnn_size,
8                                           initializer=tf.random_uniform_initializer(-0.1, 0.1, seed=2))
9             dec_cell = tf.contrib.rnn.DropoutWrapper(lstm,
10                                                    input_keep_prob = keep_prob)
11
12            # projection_layer
13            output_layer = Dense(units=vocab_size,
14                                activation='linear',
15                                kernel_initializer = tf.truncated_normal_initializer(mean = 0.0, stddev=0.1))
16
17            attn_mech = tf.contrib.seq2seq.BahdanauAttention(rnn_size,
18                                                         enc_output,
19                                                         text_length,
20                                                         normalize=False,
21                                                         name='BahdanauAttention')
22
23            dec_cell = tf.contrib.seq2seq.AttentionWrapper(dec_cell,
24                                                         attn_mech,
25                                                         rnn_size)
26
27            initial_state = dec_cell.zero_state(dtype=tf.float32, batch_size=batch_size)
28
29            with tf.variable_scope("decode"):
30                training_logits = training_decoding_layer(dec_embed_input,
31                                                         summary_length,
32                                                         dec_cell,
33                                                         initial_state,
34                                                         output_layer,
35                                                         vocab_size,
36                                                         max_summary_length)
37
38                with tf.variable_scope("decode", reuse=True):
39                    inference_logits = inference_decoding_layer(embeddings,
40                                                                vocab_to_int['<GO>'],
41                                                                vocab_to_int['<EOS>'],
42                                                                dec_cell,
43                                                                initial_state,
44                                                                output_layer,
45                                                                max_summary_length,
46                                                                batch_size)
47
48            return training_logits, inference_logits

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

1 def seq2seq_model(input_data, target_data, keep_prob, text_length, summary_length, max_summary_length,
2                   vocab_size, rnn_size, num_layers, vocab_to_int, batch_size):
3     """Use the previous functions to create the training and inference logits"""
4
5     # Use Numberbatch's embeddings and the newly created ones as our embeddings
6     embeddings = word_embedding_matrix
7
8     enc_embed_input = tf.nn.embedding_lookup(embeddings, input_data)
9     enc_output, enc_state = encoding_layer(rnn_size, text_length, num_layers, enc_embed_input, keep_prob)
10
11     dec_input = process_encoding_input(target_data, vocab_to_int, batch_size)
12     dec_embed_input = tf.nn.embedding_lookup(embeddings, dec_input)
13
14     training_logits, inference_logits = decoding_layer(dec_embed_input,
15                                                       embeddings,
16                                                       enc_output,
17                                                       enc_state,
18                                                       vocab_size,
19                                                       text_length,
20                                                       summary_length,
21                                                       max_summary_length,
22                                                       rnn_size,
23                                                       vocab_to_int,
24                                                       keep_prob,
25                                                       batch_size,
26                                                       num_layers)
27
28     return training_logits, inference_logits

```

```

1 def find_max(data):
2     return max([len(d) for d in data])
3
4
5 def pad_sentence_batch(sentence_batch, max_sentence=None):
6     """Pad sentences with <PAD> so that each sentence of a batch has the same length"""
7     if max_sentence is None:
8         max_sentence = find_max(sentence_batch)
9     else:
10        pass
11
12    sent_pad = []
13    for sentence in sentence_batch:
14        sentence = sentence + [vocab_to_int['<PAD>']] * (max_sentence - len(sentence))
15        sentence = sentence[:max_sentence]
16        sent_pad.append(sentence)
17    return sent_pad

```

```

1 def get_batches(summaries, texts, batch_size):
2     """Batch summaries, texts, and the lengths of their sentences together"""
3     len_ = len(texts) // batch_size
4     if len_ is 0:
5         raise "Error: dataset size < batch size, performance of model to low"
6
7     for batch_i in range(0, len_):
8         start_i = batch_i * batch_size
9         summaries_batch = summaries[start_i:start_i + batch_size]
10        texts_batch = texts[start_i:start_i + batch_size]
11        pad_summaries_batch = np.array(pad_sentence_batch(summaries_batch))
12        pad_texts_batch = np.array(pad_sentence_batch(texts_batch))
13
14        # Need the lengths for the _lengths parameters
15        pad_summaries_lengths = []
16        for summary in pad_summaries_batch:
17            pad_summaries_lengths.append(len(summary))
18
19        pad_texts_lengths = []
20        for text in pad_texts_batch:
21            pad_texts_lengths.append(len(text))
22
23        yield pad_summaries_batch, pad_texts_batch, pad_summaries_lengths, pad_texts_lengths

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Hyperparameters

```

1 # Set the Hyperparameters
2 batch_size = 256
3 rnn_size = 256
4 num_layers = 2
5 learning_rate = 0.005
6 epochs = 500
7
8 # Train the Model
9 learning_rate_decay = 0.995
10 min_learning_rate = 0.0001
11 display_step = 20 # Check training loss after every 20 batches
12 per_epoch = 3 # Make 3 update checks per epoch
13 update_check = (len(X_train)//batch_size//per_epoch)-1
14 stopwatch = 0
15 stop_early = 0
16 update_loss = 0
17 batch_loss = 0
18 summary_update_loss = [] # Record the update losses for saving improvements in the model

```

dropout

[x] 0%
 [] 20%
 [] 35%

optimizer

[] RMSProp
 [x] Adam

```

1 keep_probability = 1.0 # dropout
2 optimizer = tf.train.AdamOptimizer(learning_rate) # optimizer

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่าจะกรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

1 # Build the graph
2 train_graph = tf.Graph()
3 # Set the graph to default to ensure that it is ready for training
4 with train_graph.as_default():
5
6     # Load the model inputs
7     input_data, targets, lr, keep_prob, summary_length, max_summary_length, text_length = model_inputs()
8
9     # Create the training and inference logits
10    training_logits, inference_logits = seq2seq_model(tf.reverse(input_data, [-1]),
11                                                    targets,
12                                                    keep_prob,
13                                                    text_length,
14                                                    summary_length,
15                                                    max_summary_length,
16                                                    len(vocab_to_int)+1,
17                                                    rnn_size,
18                                                    num_layers,
19                                                    vocab_to_int,
20                                                    batch_size)
21
22    # Create tensors for the training logits and inference logits
23    training_logits = tf.identity(training_logits.rnn_output, 'logits')
24    inference_logits = tf.identity(inference_logits.sample_id, name='predictions')
25
26    # Create the weights for sequence_loss
27    masks = tf.sequence_mask(summary_length, max_summary_length, dtype=tf.float32, name='masks')
28
29    with tf.name_scope("optimization"):
30        # Loss function
31        # the last layer is FC NET as vocab_size units,
32        # and use activation as softmax function
33        cost = tf.contrib.seq2seq.sequence_loss(
34            training_logits,
35            targets,
36            masks)
37
38        # Gradient Clipping
39        gradients = optimizer.compute_gradients(cost)
40        capped_gradients = [(tf.clip_by_value(grad, -5., 5.), var) for grad, var in gradients if grad is not None]
41        train_op = optimizer.apply_gradients(capped_gradients)
42
43    print("Graph is built.")

```

.เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Train the Model

```

1 checkpoint = "best_model.ckpt"
2 with tf.Session(graph=train_graph) as sess:
3     sess.run(tf.global_variables_initializer())
4
5     for epoch_i in range(1, epochs+1):
6         print('=' * 70)
7         print('Epoch {:>3}/{}'.format(epoch_i, epochs))
8         print('=' * 70)
9
10        update_loss = 0
11        batch_loss = 0
12        for batch_i, (summaries_batch, texts_batch, summaries_lengths, texts_lengths) in enumerate(
13            get_batches(y_train, X_train, batch_size)):
14            start_time = time.time()
15
16            _, loss = sess.run([train_op, cost],
17                               {input_data: texts_batch,
18                                targets: summaries_batch,
19                                lr: learning_rate,
20                                summary_length: summaries_lengths,
21                                text_length: texts_lengths,
22                                keep_prob: keep_probability})
23
24            batch_loss += loss
25            update_loss += loss
26            end_time = time.time()
27            batch_time = end_time - start_time
28            stopwatch += batch_time
29
30            if (batch_i % display_step == 0) and (batch_i is not 0):
31                print('Epoch {:>3}/{} Batch {:>4}/{} - Loss: {:>6.3f}, Seconds: {:>4.2f}'
32                      .format(epoch_i,
33                               epochs,
34                               batch_i,
35                               len(X_train) // batch_size,
36                               batch_loss / display_step,
37                               batch_time * display_step))
38                batch_loss = 0
39
40            if batch_i % update_check == 0:
41                print("Average loss for this update:", round(update_loss/update_check,3))
42                summary_update_loss.append(update_loss)
43
44                # If the update loss is at a new minimum, save the model
45                if update_loss <= min(summary_update_loss):
46                    print('\nNew Record!')
47                    stop_early = 0
48                    saver = tf.train.Saver()
49                    saver.save(sess, checkpoint)
50
51                else:
52                    stop_early += 1
53                    print("No Improvement. " + str(stop_early))
54                    update_loss = 0
55
56                # Reduce learning rate, but not below its minimum value
57                learning_rate *= learning_rate_decay
58                if learning_rate < min_learning_rate:
59                    learning_rate = min_learning_rate
60
61
62        print('\n--- Finish ---')
63        print(f'Stopwatch: {stopwatch:.2f} sec.')

```

Loss Plot

```
1 import matplotlib.pyplot as plt
```

```
1 num_update_check_per_epoch = (len(X_train) // batch_size // update_check) + 1 # step update check
2 step_add = 1 / num_update_check_per_epoch
3 len_data_plot = len(summary_update_loss)
4 len_data_x_axis = len_data_plot * step_add
5 data_x_axis = np.arange(step_add, len_data_x_axis+step_add, step_add)
```

```
1 summary_update_loss = np.array(summary_update_loss) / update_check
```

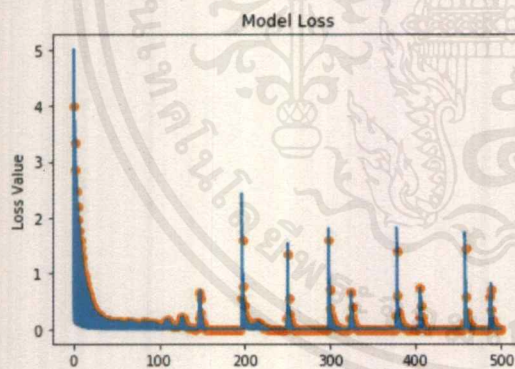
```
1 highlight_x = []
2 highlight_y = []
3 for i, (x, y) in enumerate(zip(data_x_axis, summary_update_loss)):
4     if (i+1) % 4 == 0 and x > 0:
5         highlight_x.append(x)
6         highlight_y.append(y)
```

```
1 print('number update check every epoch: ', num_update_check_per_epoch)
2 print('update check: ', update_check)
```

number update check every epoch: 4
update check: 26

Full plot

```
1 # point
2 plt.plot(highlight_x, highlight_y, 'bo', color='C1')
3
4 # line
5 plt.plot(data_x_axis, summary_update_loss, 'k', color='C0')
6
7 plt.title('Model Loss')
8 plt.ylabel('Loss Value')
9 plt.xlabel('Epoch')
10
11 plt.savefig('model_loss_full.png') # save to image
12 plt.show()
```



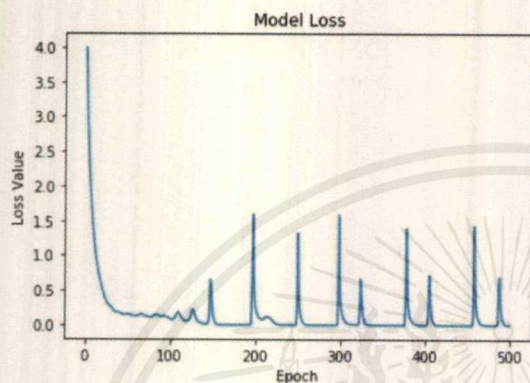
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Last loss on epoch (avg)

```

1 loss_per_epoch = [summary_update_loss[i-1] for i in range(1,
2                     len(summary_update_loss)+1) if i % num_update_check_per_epoch is 0]
3 data_x_axis = list(range(1, len(loss_per_epoch)+1))
4
5 # line
6 plt.plot(data_x_axis, loss_per_epoch, 'k', color='C0')
7
8 plt.title('Model Loss')
9 plt.ylabel('Loss Value')
10 plt.xlabel('Epoch')
11
12 plt.savefig('model_loss_epoch.png') # save to image
13 plt.show()

```



```

1 to_txt('summary_update_loss.txt', summary_update_loss)

```

Model Evaluate

```

1 def text_to_seq(text):
2     """Prepare the text for the model"""
3     return [vocab_to_int.get(word, vocab_to_int['<UNK>']) for word in text]

```

```

1 # set generate summaries length
2 generate_len = find_max(y_train)
3
4 print('# Train set')
5 print("Text length: ", find_max(X_train))
6 print("Summaries length: ", generate_len)
7 print('-----')
8 print('# Test set')
9 print("Text length: ", find_max(X_test))
10 print("Summaries length: ", find_max(y_test))

```

```

# Train set
Text length: 49
Summaries length: 26

```

```

-----
# Test set
Text length: 43
Summaries length: 22

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

1 texts = X_test
2 checkpoint = "./best_model.ckpt"
3
4 loaded_graph = tf.Graph()
5 with tf.Session(graph=loaded_graph) as sess:
6     # Load saved model
7     loader = tf.train.import_meta_graph(checkpoint + '.meta')
8     loader.restore(sess, checkpoint)
9
10    input_data = loaded_graph.get_tensor_by_name('input:0')
11    logits = loaded_graph.get_tensor_by_name('predictions:0')
12    text_length = loaded_graph.get_tensor_by_name('text_length:0')
13    summary_length = loaded_graph.get_tensor_by_name('summary_length:0')
14    keep_prob = loaded_graph.get_tensor_by_name('keep_prob:0')
15
16    outputs = []
17    for text in texts:
18        #Multiply by batch_size to match the model's input parameters
19        answer_logits = sess.run(logits,
20                                {input_data: [text]*batch_size,
21                                 summary_length: [generate_len],
22                                 text_length: [len(text)]*batch_size,
23                                 keep_prob: keep_probability})[0]
24        outputs.append(answer_logits)

```

```
1 print('Number of output: ', len(outputs))
```

Number of output: 2377

- **ROUGE** : metrics for Text summarization

```

1 import sys
2 sys.path.append('../..')
3
4 from evaluate import rouge
5
6 def clean_pad(lists, pad_int=vocab_to_int['<PAD>']):
7     result = []
8     for list_ in lists:
9         without_pad = [l for l in list_ if l != pad_int]
10        result.append(without_pad)
11    return result

```

```

1 # remove padding in outputs (generate) for ROUGE
2 outputs_without_pad = clean_pad(outputs)
3
4 # covert to vocab
5 y_test_text = [[int_to_vocab[i] for i in y] for y in y_test]
6 outputs_text = [[int_to_vocab[i] for i in output] for output in outputs_without_pad]
7
8 # join vocab to sentent separate by '|' charactor
9 y_test_text = ["|".join(text) for text in y_test_text]
10 outputs_text = ["|".join(text) for text in outputs_text]

```

```

1 # compute ROUGE score
2 report = rouge(hypotheses=outputs_text, references=y_test_text)
3 print(report)

```

```

{'rouge_1/f_score': 0.22370522924266248,
'rouge_1/r_score': 0.21714887438573277,
'rouge_1/p_score': 0.245603127273964,
'rouge_2/f_score': 0.07147199573550303,
'rouge_2/r_score': 0.07299463691164472,
'rouge_2/p_score': 0.07520423450945052,
'rouge_l/f_score': 0.1854591933014478,
'rouge_l/r_score': 0.20306315305196806,
'rouge_l/p_score': 0.200269883421336}

```

```
1 to_json('report.json', report)
```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Summary Result

```

1 # covert to vocab
2 X_test_text = [[int_to_vocab[i] for i in x] for x in X_test]
3
4 # join vocab to sentent separate by ' ' charactor, without <EOS>
5 X_test_text = [" ".join(text[:-1]) for text in X_test_text]

```

```

1 for text, summ_true, summ_gen in zip(X_test_text, y_test_text, outputs_text):
2     print("Text ---> ', text)
3     print("True summary ---> ', summ_true)
4     print("Generate summary ---> ', summ_gen)
5     print('-' * 70)

```

Text ---> ผู้บัญชาการตำรวจแห่งชาติ|خانรับ|คำสั่ง|นายกรัฐมนตรี|ห้าม|ผู้ต้องหา|แถลงข่าว|สื่อมวลชน|ชี้|ปฏิบัติตาม|พิจารณา|ข้อ
บกพร่อง|วินัย|อาญา

True summary ---> ผบ.ตร.|รับ|ลูก|นายกฯ|ห้าม|นำ|ผู้ต้องหา|มา|นั่ง|แถลงข่าว

Generate summary ---> ผบ.ตร.|สั่ง|ให้|นำ|เกสโอดู|คำสั่ง|คสช.|ผิดวินัย|หลัง|ถูก|ตัดสิน

Text ---> ทช.|เดือน|แมงกะพรุน|กลอง|จำนวนมาก|บริเวณ|ชายหาด|เกาะ|หมาก|จ.|ตราด|นักท่องเที่ยว|สัมผัส|โดน|บริเวณ|เขา|ปร
มพยาบาล|เบื้องต้น|ชี้|แมงกะพรุน|พิษ|ปฐมพยาบาล|วิธี|เสี่ยง|ตาย|แนะ|ผู้ประกอบการ|แจ้ง|เดือน|นักท่องเที่ยว|เตรียม|นำ|สนสายช|ให้
าม

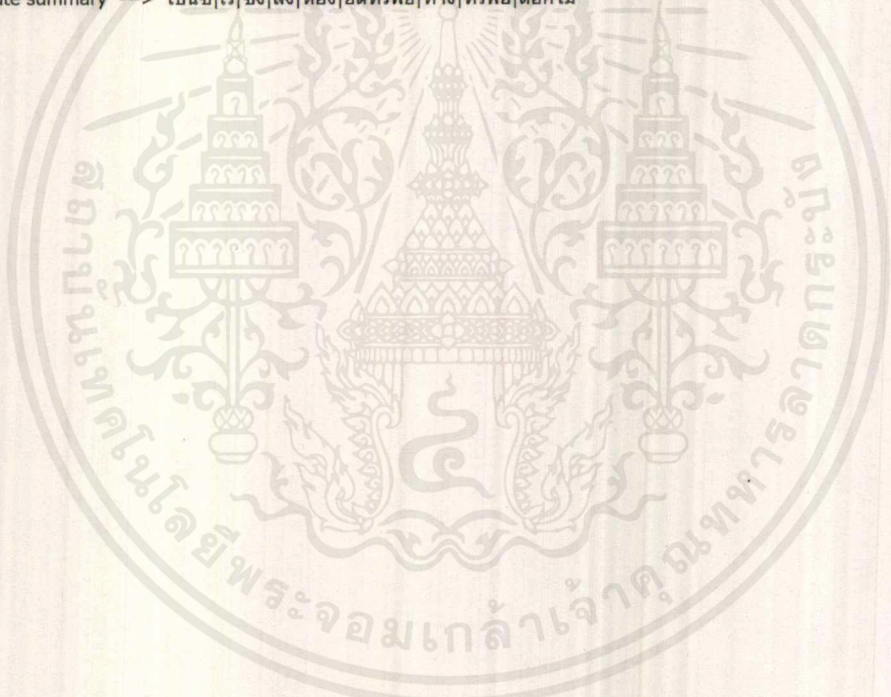
True summary ---> แมงกะพรุน|กลอง|โผล่|เกาะ|หมาก|เดือน|พิษ|ถึงตาย

Generate summary ---> เดือน|เล่น|นำ|หาด|เกาะ|ก้อน|ห้วง|ชนเขา|โดน|พิษ

Text ---> อัศจรรย์|วัด|โรจน์|เจริญ|เดช|เบนซ์|เร|ซึ่ง|มารดา|หลักฐาน|การครอบครอง|รถยนต์|ทรู|ทรัพย์สิน|ตำรวจ|ปราบปราม|ยา
เสพติด|ตรวจสอบ

True summary ---> เบนซ์|เร|ซึ่ง|หอบ|หลักฐาน|ครอบครอง|รถ|ทรู|ให้|ตำรวจ|ปส.

Generate summary ---> เบนซ์|เร|ซึ่ง|ส่ง|ห้อง|ยึดทรัพย์สิน|ทาง|ทรัพย์สิน|ดอกไม้



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

โค้ดสร้างแบบจำลองประเภทอื่น ๆ มีความเหมือนกันกับโค้ดสร้างแบบจำลอง Adam without Dropout Model 1 ซึ่งจะมีความแตกต่างเพียงการกำหนดค่า Hyperparameter ที่อยู่ในส่วนของการทำรอบสีแดงของโค้ดดังกล่าว

2. โค้ดสร้างแบบจำลอง Adam without Dropout Model 2 เหมือนกันกับโค้ดสร้างแบบจำลอง Adam without Dropout Model 1

3. โค้ดกำหนดค่า Hyperparameter ของ Adam Dropout 0.2

dropout

0%
 20%
 35%

optimizer

RMSProp
 Adam

```
1 keep_probability = 0.80 # dropout
2 optimizer = tf.train.AdamOptimizer(learning_rate)
```

4. โค้ดกำหนดค่า Hyperparameter ของ Adam Dropout 0.35

dropout

0%
 20%
 35%

optimizer

RMSProp
 Adam

```
1 keep_probability = 0.65 # dropout
2 optimizer = tf.train.AdamOptimizer(learning_rate) # optimizer
```

5. โค้ดกำหนดค่า Hyperparameter ของ RMSProp without Dropout

dropout

0%
 20%
 35%

optimizer

RMSProp
 Adam

```
1 keep_probability = 1.0 # dropout
2 optimizer = tf.train.RMSPropOptimizer(learning_rate)
```

6. ได้กำหนดค่า Hyperparameter ของ RMSProp Dropout 0.2

dropout

- 0%
- 20%
- 35%

optimizer

- RMSProp
- Adam

```
1 keep_probability = 0.80 # dropout
2 optimizer = tf.train.RMSPropOptimizer(learning_rate)
```

7. ได้กำหนดค่า Hyperparameter ของ RMSProp Dropout 0.35

dropout

- 0%
- 20%
- 35%

optimizer

- RMSProp
- Adam

```
1 keep_probability = 0.65 # dropout
2 optimizer = tf.train.RMSPropOptimizer(learning_rate)
```

8. โค้ดสร้าง Dataset

```

1 from nltk.corpus import words as nltk_words
2 from nltk.corpus import stopwords
3 import pythainlp
4 from pythainlp import tokenize
5 from pythainlp.corpus import thai_stopwords
6 from num2words import num2words
7 import matplotlib.pyplot as plt
8 import pandas as pd
9 import numpy as np
10 import re
11 import time
12
13 print('PyThaiNLP Version: {}'.format(pythainlp.__version__))

```

PyThaiNLP Version: 2.0.5

Inspecting the Data

```

1 dataset_path = './dataset/thaipbs_news/scraping/news_list/thaipbs_news_dataset.csv'
2 language = 'thai'
3 tk_engine = 'newmm'

```

```

1 # dtype: dict
2 custom_contractions_dict = {
3     '%': 'เปอร์เซ็นต์',
4     '\xa0': ' ',
5 }
6
7 # dtype: set
8 custom_dictionary = {
9     '"',
10    "'",
11    'ซีพี',
12    'คิวดี',
13    'อีอีซี',
14    'มธ.',
15    'แซชแท็ก',
16    'ส.ส.ท.',
17    'อีไอดี',
18    'ฟลุ่มน',
19    'ไอช',
20    'ริเมค',
21    'แมงค็อก',
22    'เอฟซี',
23    'ทีซีดีซี',
24    'เฟชมค',
25    'สคร.',
26    'นูโร',
27    'กยศ.',
28 }
29
30 # dtype: set
31 custom_stopwords = {
32     'รก',
33 }

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```
1 df = pd.read_csv(dataset_path)
```

```
1 display(df.shape)
2 display(df.head())
```

(29170, 2)

	highlights	story
0	ส่งข้อมูล WhatsApp โยงหากรรมกรนักข่าวฯฯ	ผู้ใช้งานแอปพลิเคชัน "วอทสแอฟ" ทั่วโลกประมาณ 1...
1	สถาบันศึกษา มธ.ยื่นจดหมายขอยกเลิกบังคับใส่เครี...	สถาบันศึกษามหาวิทยาลัยธรรมศาสตร์ศูนย์รังสิต ยี...
2	สรรพสามิตแจงเพิ่มภาษีดีเซล หวังหนุนใช้ป่าสมผลิ...	กรมสรรพสามิตปรับเพิ่มอัตราภาษีน้ำมันดีเซล หวัง...
3	ท่าแผนฆ่าสามมีฝั่งดิน อ้างถูกขูฆ่า-ทำร้ายลูก	ผู้ต้องหาอ้างมีสามมีเสียชีวิตและนำศพไปฝังดิน ท...
4	"ไดโซเดียม" ในอาหารกำลังเร่จรูไปโซผงชูรส	แพทย์ยืนยัน ไดโซเดียม ส่วนประกอบในขนมปังและโจ้...

```
1 # Check for any nulls values
2 df.isnull().sum()
```

```
highlights 0
story 0
dtype: int64
```

```
1 from sklearn.utils import shuffle
```

```
1 df.columns = ['Summary', 'Text']
2
3 # Remove null values and unneeded features
4 df = df.drop_duplicates().dropna() # drop nan & duplicates
5 df = shuffle(df) # shuffle data
6 df = df.reset_index(drop=True)
```

```
1 df.head()
```

	Summary	Text
0	เงินค่านสหรัฐฯ ขยายอาวาร์ไฟใต้หวัน	เงินเดือนสหรัฐฯ ว่าความล้มพินธ์ทวิภาค้อาจสั้นคล...
1	"พล.ต.อ.สมยศ-โคร์แซง-ซีโก" รวมประมุขเตรียมที่...	สมาคมกีฬาฟุตบอลวางแผนเตรียมทีมฟุตบอลทีมชาติไท...
2	อนาคต "วัดพระธรรมกาย" หากไร "พระธัมมชโย"	อนาคต "วัดพระธรรมกาย" หากไม่มี "พระธัมมชโย" ศร...
3	"จิสต้า" ไซค์ดาวเทียม 10 ดวงติดตามไฟป่าตั้งเป้า...	กระทรวงทรัพยากรธรรมชาติและสิ่งแวดล้อม จัด2 ภาด...
4	เลือกตั้ง2562: ทอ.แจงปมเบรกพรรคอนาคตใหม่หาเสีย...	โฆษกกองทัพอากาศ ยืนยันไม่ได้กีดกันนายนิคม แสง...

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Sample Data

```

1 # Inspecting some of the news
2 for i in range(5):
3     print("Review #",i+1)
4     print(df.Summary[i])
5     print(df.Text[i])
6     print()

```

Review # 1

จีนค้าสหรัฐฯ ขายอาวุธให้ไต้หวัน

จีนเดือนสหรัฐฯ ว่าความสัมพันธ์ทวิภาคีอาจสั่นคลอนหากสหรัฐฯ ขายอาวุธให้ไต้หวัน พร้อมย้ำถึงนโยบายจีนเดียวและจะไม่ยอมให้มีใครมาแยกดินแดนอย่างเด็ดขาด

Review # 2

"พล.ต.อ.สมยศ-โค้ชเฮง-ซีโก" ร่วมประชุมเตรียมทีมลุยศึกฟุตบอลโลก รอบคัดเลือก 12 ทีมสุดท้ายโซนเอเชีย

สมาคมกีฬาฟุตบอลฯวางแผนเตรียมทีมฟุตบอลทีมชาติไทยก่อนลุยศึกฟุตบอลโลก รอบคัดเลือก 12 ทีมสุดท้ายโซนเอเชีย

Review # 3

อนาคต "วัดพระธรรมกาย" หากไร้ "พระธัมมชโย"

อนาคต "วัดพระธรรมกาย" หากไม่มี "พระธัมมชโย" ศรีธธายัง "แข็งแกร่ง" หรือ "สั่นคลอน" สำนักงานพระพุทธศาสนาแห่งชาติ เชื้อผลที่เก็ดขึ้นกับพระธัมมชโยไม่กระทบกับกิจกรรมของวัด ขณะที่อดีตลูกศิษย์เชื้อสถานภาพของวัดอาจถูกสั่นคลอน โดยเฉพาะการบริหารจัดการภายในวัด

Review # 4

"จีสด้า" ไข่ดาวเทียม 10 ดวงติดตามไฟฟ้าตั้งเป้าลดร้อยละ 30

กระทรวงทรัพยากรธรรมชาติและสิ่งแวดล้อม จัด 2 มาตรการส่งเสริมประมงมือเผาป่าต้นเห็ดหมอกควินพีชในภาคเหนือ พร้อมคาดโทษจนท.รณผิดชอบในการลดจุดความร้อนระดับพื้นที่ให้ได้ร้อยละ 30 เดือน.ค.นี้ถ 5 ขาดแก้มหมอกควินข้ามแดน

Review # 5

เลือกตั้ง 2562: ทอ.แจงบมเบรกพรรคอนาคตใหม่หาเสียงเขต รร.

โฆษกกองทัพอากาศ ยืนยันไม่ได้กีดกันนายนิกร แสงศิริวานัน ผู้สมัคร ส.ส.กรุงเทพมหานคร พรรคอนาคตใหม่ ที่ออกมาระบุว่าถูกสว.รทรทหารกองทัพอากาศ (สท.) ไม่ให้เข้าไปหาเสียงในโรงเรียนฤทธิยะวรรณาลัย ซ้อยู่ในพื้นที่กองทัพอากาศ หากหาเสียงต้องขออนุญาต

Preparing the Data

```

1 import json
2
3 def to_json(fname, data):
4     with open(fname, 'w') as file:
5         json_data = json.dumps(data)
6         file.write(json_data)
7
8 def read_json(fname):
9     with open(fname, 'r') as file:
10        data = file.read()
11        json_data = json.loads(data)
12        return json_data
13
14 def to_txt(fname, data):
15     with open(fname, 'w', encoding='utf-8') as file:
16        data = [str(d) for d in data]
17        file.write('\n'.join(data))
18
19 def read_txt(fname):
20     with open(fname, 'r', encoding='utf-8') as file:
21        data = file.read()
22        string_data = data.split('\n')
23        try:
24            float_data = [float(s) for s in string_data]
25            return float_data
26        except ValueError:
27            return string_data

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

1 # A list of contractions from:
2 # http://stackoverflow.com/questions/19790188/expanding-english-language-contractions-in-python
3
4 contractions_dict = read_json('../contractions.json')
5 contractions_dict.update(custom_contractions_dict)
6
7
8 def contractions(text):
9     # replace contractions with their longer forms
10    for word1, word2 in contractions_dict.items():
11        text = text.replace(word1, word2)
12    return text

```

```

1 contractions_dict.update(custom_contractions_dict)
2 th_dictionary = set(tokenize.DEFAULT_DICT_TRIE) # load default dictionary
3 english_dictionary = set(nltk_words.words())
4 DICTIONARY = th_dictionary.union(english_dictionary).union(custom_dictionary)
5 TK = tokenize.Tokenizer(DICTIONARY, engine=tk_engine)

```

```

1 # stopwords
2 english_stopwords = set(stopwords.words('english'))
3 thai_stopwords_ = set(thai_stopwords())
4 stopwords_ = english_stopwords.union(thai_stopwords_).union(custom_stopwords)

```

```

1 to_txt('stopwords.txt', stopwords_)
2 to_txt('dictionary.txt', DICTIONARY)
3 to_json('contractions_dict.json', contractions_dict)

```

```

1 def clean_token(tokens, stopword_dict={}, use_remove_symbol=True):
2     clean = []
3     for token in tokens:
4         token = token.strip(' ') # remove white space
5         if token in stopword_dict:
6             continue
7
8         if use_remove_symbol:
9             token = remove_symbol(token)
10        if len(token) <= 1:
11            continue
12
13        clean.append(token)
14    return clean

```

```

1 def build_vocab(text, unk_method='replace', use_norm_num=True, **kwargs):
2     """
3     unk_method: unknow vocab method
4     replace: replace know vocab to <unk>
5     remove: pass vocab (remove vocab)
6     ignore: can not operate
7     """
8
9     # to lower-case
10    text = text.lower()
11    # replace contractions with their longer forms
12    text = contractions(text)
13    # word tokenize
14    tokens = TK.word_tokenize(text)
15
16    segments = []
17    for token in tokens:
18        if use_norm_num:
19            # convert number to word
20            token_list = norm_num(token, **kwargs)
21        else:
22            token_list = [token]
23        # append to segments list
24        segments += token_list
25    return segments

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

1 def norm_num(string, use_num2word=False, lang='thai'):
2     "return dtype: list"
3
4     ignore_string = ['infinity']
5
6     if string in ignore_string:
7         return string
8
9     try:
10        new_string = string.replace(',', '')
11        num = float(new_string)
12        try:
13            num = int(new_string)
14        except ValueError:
15            pass
16
17        if use_num2word:
18            word = num2words(num, lang=lang) # to word
19            tokens = TK.word_tokenize(word)
20        else:
21            tokens = [str(num)]
22        return tokens # dtype: list
23    except ValueError:
24        return [string]

```

```

1 def remove_symbol(word):
2     # Format words and remove unwanted characters
3     word = re.sub(r'https?:\V.*[\r\n]*', "", word, flags=re.MULTILINE)
4     word = re.sub(r'\<a href', "", word)
5     word = re.sub(r'&#', "", word)
6     # word = re.sub(r'[_"-;%()]+&=*,!?:#$@|[\V/]', "", word) # with . symbol
7     word = re.sub(r'[_"-;%()]+&=*,!?:#$@|[\V/]', "", word) # without . symbol
8     word = re.sub(r'<br />', "", word)
9     word = re.sub(r'\V', "", word)
10    return word

```

We will remove the **stopwords from the texts** because they do not provide much use for training our model. However, we will **keep them for our summaries** so that they sound more like natural phrases.

```

1 clean_texts = [build_vocab(text, lang=language) for text in df.Text]
2
3 # clean token with stopword
4 clean_texts = [clean_token(tokens, stopword_dict=stopwords_) for tokens in clean_texts]
5 print("Texts are complete.")
6
7 clean_summaries = [build_vocab(text, lang=language) for text in df.Summary]
8
9 # clean token without stopword
10 clean_summaries = [clean_token(tokens, stopword_dict={}) for tokens in clean_summaries]
11 print("Summaries are complete.")

```

Texts are complete.
Summaries are complete.

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

1 # Inspect the cleaned summaries and texts to ensure they have been cleaned well
2 for i in range(5):
3     print("Build & Clean Review #", i+1)
4     print(clean_summaries[i])
5     print(clean_texts[i], end='\n\n')

```

Build & Clean Review # 1

['จีน', 'คำนำ', 'สหรัฐอเมริกา', 'ชาย', 'อาวุธ', 'ไฟ', 'ไต้หวัน']

['จีน', 'เดือน', 'สหรัฐอเมริกา', 'ความสัมพันธ์', 'ทวิภาคี', 'สันคอลลอน', 'สหรัฐอเมริกา', 'ชาย', 'อาวุธ', 'ไต้หวัน', 'ยา', 'นโยบาย', 'จีน', 'ยอมให้', 'แยก', 'ดินแดน', 'อย่างเด็ดขาด']

Build & Clean Review # 2

['พล.ต.อ.', 'สม', 'ยศ', 'โคชัย', 'เฮง', 'ช', 'โก', 'ร่วม', 'ประชุม', 'เตรียม', 'ทีม', 'ลุย', 'ศึก', 'ฟุตบอล', 'โลก', 'รอบ', 'คัดเลือก', '12', 'ทีม', 'สุดท้าย', 'โซน', 'เอเชีย']

['สมาคม', 'กีฬา', 'ฟุตบอล', 'วางแผน', 'เตรียม', 'ทีม', 'ฟุตบอล', 'ทีม', 'ชาติ', 'ไทย', 'ลุย', 'ศึก', 'ฟุตบอล', 'โลก', 'รอบ', 'คัดเลือก', '12', 'ทีม', 'สุดท้าย', 'โซน', 'เอเชีย']

Build & Clean Review # 3

['อนาคต', 'วัด', 'พระ', 'ธรรมกาย', 'หาก', 'ไว้', 'พระ', 'ธัมม', 'ชโย']

['อนาคต', 'วัด', 'พระ', 'ธรรมกาย', 'ไม่มี', 'พระ', 'ธัมม', 'ชโย', 'ศรัทธา', 'แข็งแกร่ง', 'สันคอลลอน', 'สำนักงาน', 'พระพุทธศาสนา', 'แห่งชาติ', 'ที่เกิด', 'ขึ้นกับ', 'พระ', 'ธัมม', 'ชโย', 'กระทบ', 'กิจกรรม', 'ของวัด', 'ลูกศิษย์', 'สถานภาพ', 'ของวัด', 'สันคอลลอน', 'โดยเฉพาะ', 'การบริหาร', 'วัด']

Build & Clean Review # 4

['จี', 'สดำ', 'ไซ', 'ดาวเทียม', '10', 'ดวง', 'ติดตาม', 'ไฟฟ้า', 'ตั้งเป้า', 'ลด', 'ร้อยละ', '30']

['กระทรวง', 'ทรัพยากรธรรมชาติ', 'สิ่งแวดล้อม', 'จัด', 'มาตรการ', 'สังคม', 'ประ', 'นาม', 'มือ', 'เผา', 'ป่า', 'ต้นเหตุ', 'หมอก', 'ควันพิษ', 'ภาคเหนือ', 'คาดโทษ', 'จนท.', 'รับผิดชอบ', 'ลด', 'จุด', 'ความร้อน', 'ระดับ', 'พื้นที่', 'ไฟใต้', 'ร้อยละ', '30', 'เดือน', 'ม.ค.', 'ถก', 'ชาติ', 'แก', 'หมอก', 'ควัน', 'ข้ามแดน']

Build & Clean Review # 5

['เลือกตั้ง', '2562', 'ทอ.', 'แจง', 'ปม', 'เมรก', 'พรรค', 'อนาคต', 'ใหม่', 'หาเสียง', 'เขต', 'รร.']

['โฆษก', 'กองทัพอากาศ', 'ไม่ได้', 'คัดค้าน', 'ถ', 'กม', 'แสง', 'สี', 'ริ', 'นาวิน', 'ผู้สมัคร', 'ส.ส.', 'กรุงเทพมหานคร', 'พรรค', 'อนาคต', 'ออกม', 'า', 'รม', 'สารวัตรทหาร', 'กองทัพอากาศ', 'สท.', 'เข้าไป', 'หาเสียง', 'โรงเรียน', 'ฤทธิยะวรรณาลัย', 'ซี', 'พื้นที่', 'กองทัพอากาศ', 'หาเสียง', 'ขออนุญาต']

```

1 def count_words(count_dict, text):
2     """Count the number of occurrences of each word in a set of text"""
3     for sentence in text:
4         for word in sentence:
5             if word not in count_dict:
6                 count_dict[word] = 1 # new key-case
7             else:
8                 count_dict[word] += 1

```

```

1 # Find the number of times each word was used and the size of the vocabulary
2 word_counts = {}
3
4 count_words(word_counts, clean_summaries)
5 count_words(word_counts, clean_texts)
6
7 print("Size of Vocabulary:", len(word_counts))

```

Size of Vocabulary: 25399

Word2Vec: PyThaiNLP

```

1 from pythainlp.word_vector import get_model
2
3 ww_model = get_model()
4 embeddings_index = {vocab: ww_model[vocab] for vocab in ww_model.vocab} # ww_dict
5
6 print("Word embeddings:", len(embeddings_index))

```

Word embeddings: 51358

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

1 # Find the number of words that are missing from CN, and are used more than our threshold.
2 missing_words = 0
3 threshold = 20
4
5 for word, count in word_counts.items():
6     if word not in embeddings_index:
7         missing_words += 1
8
9 missing_ratio = round(missing_words/len(word_counts),4)*100
10
11 print('Word counts: ', len(word_counts))
12 print("Number of words found: ", len(word_counts) - missing_words)
13 print("Number of words missing:", missing_words)
14 print("Percent of words that are missing from vocabulary: {}".format(missing_ratio))

```

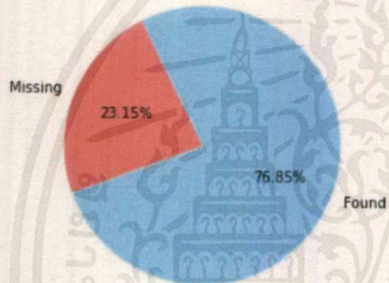
Word counts: 25399
 Number of words found: 19518
 Number of words missing: 5881
 Percent of words that are missing from vocabulary: 23.150000000000002%

```

1 labels = ['Found', 'Missing']
2 sizes = [len(word_counts) - missing_words, missing_words]
3 plt.title('Percentage of words found from vector dictionary')
4 plt.pie(sizes, labels=labels, autopct='%1.2f%%', startangle=200, colors=['lightskyblue', 'lightcoral'])
5 plt.axis('equal')
6 plt.show()

```

Percentage of words found from vector dictionary



vocab2int, int2vocab

```

1 #dictionary to convert words to integers
2 vocab_to_int = {}
3
4 value = 0
5 for word, count in word_counts.items():
6     if count >= threshold or word in embeddings_index:
7         vocab_to_int[word] = value
8         value += 1
9
10 # Special tokens that will be added to our vocab
11 codes = ["<UNK>", "<PAD>", "<EOS>", "<GO>"]
12
13 # Add codes to vocab
14 for code in codes:
15     vocab_to_int[code] = len(vocab_to_int)
16
17 # Dictionary to convert integers to words
18 int_to_vocab = {}
19 for word, value in vocab_to_int.items():
20     int_to_vocab[value] = word

```

```

1 # Need to use 300 for embedding dimensions
2 embedding_dim = 300
3 nb_words = len(vocab_to_int)
4
5 # Create matrix with default values of zero
6 word_embedding_matrix = np.zeros((nb_words, embedding_dim), dtype=np.float32)
7 new_word_embedding = []
8
9 for word, i in vocab_to_int.items():
10     if word in embeddings_index:
11         word_embedding_matrix[i] = embeddings_index[word]
12     else:
13         # If word not in thai word2vec vector dict, create a random embedding for it
14         new_embedding = np.array(np.random.uniform(-1.0, 1.0, embedding_dim))
15         embeddings_index[word] = new_embedding
16         word_embedding_matrix[i] = new_embedding
17
18     new_word_embedding.append(word)
19
20 # Check if value matches len(vocab_to_int)
21 print('number of word embedding matrix: ', len(word_embedding_matrix))
22
23 print('number of new word embedding: ', len(new_word_embedding))
24 print('new word embedding: ', new_word_embedding)

```

number of word embedding matrix: 19878

number of new word embedding: 360

new word embedding: ['ท.', 'ยครว', 'ขุนน้ำ', 'นายภา', 'สหรัฐ', 'ใบ', 'ดอนเมือง', 'ดอลลาร์สหรัฐ', 'อุตุฯ', 'ผู้ถูกกล่าวหา', 'คป', 'เพชรม', 'ผู้ว่าฯ', 'รวมใจ', 'ซีพี', 'ผิดวินัย', 'ปชช.', 'ลุมพินี', 'สป', 'เข้าให้', 'บี', 'โดนัลด์', 'ทรัมป์', 'เมียนมา', 'ผู้ต้องขัง', 'กรุงเทพฯ', 'อุทยานฯ', 'โ', 'กรม', 'เน', 'ธัมม', 'สตร.', 'โรดแมป', 'สนง.', 'สค', 'รัชต์', 'ธีร', 'อินโดฯ', 'ส.', 'ลมกระโชก', 'ลัมพิน', 'เฟส', 'โรธัญญา', 'โบศ', 'สวล', 'ส.ส.ท.', 'เดามะเฟือง', 'ลิมง', 'พะเนิน', 'โปรดเกล้าฯ', 'ลอนดอน', 'บู', 'เอฟซี', 'คพ.', 'รับซื้อ', 'จฟ้าฯ', 'กรมชล', 'อินทนนท์', '40000', 'สมเด็จพระสังฆราช', 'แม่แจ่ม', 'สุโขทัย', 'จยย.', 'โตะ', 'รอมฎอน', 'ภูมิพล', 'น่านาคาล', 'โจชัว', 'หวอง', 'ดุน', 'ซูเปอร์', 'เรย์', 'วิทย์', 'ทศ.', 'ก', 'ยท.', 'หัด', 'มาร์ก', 'ธา', 'กองปราบฯ', 'กรมอุทยานฯ', 'กรุงเทพ', 'ชฎ', 'ดรัส', 'จช', 'วธ', '100000', 'ยังมี', 'ดี', 'พีดีอาร์ไอ', '12000', 'เงิน', 'ทอน', 'hia', 'อร', 'ธป', 'คอลเซ็นเตอร์', 'ตั้ง', 'บร', 'ณ', 'เบลเยียม', 'ณัฐ', 'แอดมิน', 'สระแก้ว', 'สวีเดน', 'คปท.', 'สพฉ.', 'แอฟฟลิเคชัน', 'ร', 'นมาชนบช', 'เปอร์', 'ธนา', 'เอ็นจีวี', 'ช.', 'แชมป์', '300000', 'สจ', 'แกริม', 'รับงาน', 'ทับซ้อน', 'สเปน', 'นญา', 'สปส.', 'ภูพาน', 'สมยศ', 'กร', 'ระลอก', 'ไซเบอร์', 'วีร', 'เรือโดยสาร', 'อัล', 'บร', 'สพฐ.', 'สคม.', 'คอนฯ', 'ธา', 'ท', 'ทปอ.', 'ผู้ตรวจการแผ่นดิน', 'ล่อมคอก', 'ส', 'ตำรวจ', 'างหลวง', '60000', 'ปท.', 'เลสเด', 'สังข', 'อี', 'แพทย์สภา', 'สุรเชษฐ์', 'กทพ.', 'พะยุง', 'เวบ', 'gat', 'ทช.', 'ช.ภ.', 'ธนะ', 'สดี', 'เจ้าท่า', 'เ', 'ลชชา', '13000', 'ปลอดเลือด', 'น', 'ทร', 'ตู้เอทีเอ็ม', 'ณป', 'ณ', 'หนองบัวลำภู', 'ปานามา', 'แอล', 'กยศ.', 'บท.1', 'อี', 'พงษ์', 'ชิตี', 'อิสติน', 'บูล', 'ท้ายเขื่อน', 'เซอร์เคน', 'พงษ์', 'เชิน', 'กรมบังคับคดี', 'สอนพยาน', 'เน', 'อีค', 'โรล', 'ปวีณา', 'ออง ซาน ซู จี', 'มิ่งกาฬ', 'ทวิตเตอร์', 'อบลฯ', 'วันเฉลิมพระชนมพรรษา', 'ผือ', 'เชลฟี', 'เนจิก', 'โปเกมอน', 'สท', 'โลลา', 'วีร', 'กฤษณ', 'ยอนคร', 'หัวคิว', 'ลาสเวกัส', 'ทร', 'ศ', 'ป', 'มทิลล', 'ประ', 'สก', '50000', 'ภาษาสรรพสามิต', '70000', 'สิดอ', 'ซาอุฯ', 'กิน', 'บ่งอยู่', 'เฮซ', 'บี', 'กส', 'หัย', 'นาปรัง', 'เป้าหมายเสียง', 'วังช้าง', 'ฟองกลิ่น', 'ดอยด', 'ogle', 'ไอซี', 'สว', 'วสท.', 'นัม', 'นค', 'เข', '80000', 'ข้างแอฟริกา', 'จับหน่อม', 'ธงชาติไทย', '35000', 'ร', 'ช', 'กส', 'พายุไซโคลน', 'แพทย์', 'วัฒน', '15000', 'ออกโฉนด', 'ซูจี', 'อดุลยเดช', 'สุนัขจรจัด', 'หมายค้น', 'พรากผู้เยาว', 'ให้ความเป็นธรรม', 'อี', 'โบท', 'ดดม', 'สร', 'รัตน', 'สต็อก', 'พัฒนา', 'เอกวาดอร์', 'สตรีท', 'แอฟฟลิเคชัน', 'อีพ', 'โอบามา', 'สทช.', '14.0', 'ปีศาจ', '2000', '00', 'แจ็ด', 'จตุพร', 'วี', 'เนลดี', 'แคนดี้เดด', 'โทษทางวินัย', 'กระสอบทราย', 'พญาไท', 'อีซีซี', 'สันดา', 'โบอิง', 'ทีฟ', 'ปุดิน', 'โทรโยค', 'สถานการณ์ฉุกเฉิน', 'แคมป์', 'ซิมซุง', 'ตัดประเด็น', 'ยะไซ', 'อี', 'ดีเบต', 'บรัสเซลส์', 'เกีย', 'คลื่นความถี่', 'คืนให้', 'กรมการจัดหางาน', 'วัน', 'สงกรานต์', 'โลโก้', '22.0', 'พาร', 'ซาอ', 'อาร์เมเนีย', 'ราค', 'ซัน', 'ซีบีบีไอ', 'สายไฟฟ้า', 'นาคะ', 'เป็นราย', 'โนมินี', 'พวล', 'ผู้ประกันตน', 'ชินวัตร', 'ธัญบุรี', 'เสด็จสวรรคต', 'วางมาตรการ', 'ดิญ', 'จุฬาลงกรณ์', 'hall', 'อินสตาแกรม', 'บัน', 'มาร์ค', 'ไทใต้', 'ป่าสงวนแห่งชาติ', 'ออก', 'ของ', 'ถึงชีวิต', 'สท', '500000', 'คช', 'ไรต์', 'วีช', 'ไม่ท้วงห้าม', '18.0', 'เก', 'อี', 'ส่วนราชการ', 'ปากทอ้ง', 'ทำเนิน', 'ยอดน้ำ', 'ล้าน', 'า', 'ยอดตลาด', 'รับของโจร', '600000', 'กระทรวงเกษตรฯ', 'หน้า', 'ละม่ง', 'สีลอด', '400000', 'สาหวัย', 'แปรม ดินสุสานนท์', 'การออกเสี', 'ยประชาบดี', 'อยู่แล้ว', 'ผู้ดำรงตำแหน่งทางการเมือง', 'วี', 'พระราชบัญญัติประกอบรัฐธรรมนูญ', 'เพีย', 'ทฤษฎี สหวงษ์', 'สมเด็จพระเท', 'พรัตนราชสุดาฯ สยามบรมราชกุมารี', 'พราหมณ์', 'นายสุเทพ เทือกสุบรรณ', '17.0', 'ว่าไม่ได้', 'ตลอด 24 ชั่วโมง', 'ต่าง ๆ', 'กิตติ', '15.0', 'โดยที่', '16.0', 'นายอภิสิทธิ์ เวชชาชีวะ', '<UNK>', '<PAD>', '<EOS>', '<GO>']

```

1 to_json('vocab_to_int.json', vocab_to_int)
2 to_json('int_to_vocab.json', int_to_vocab)
3 np.save('word_embedding_matrix.npy', word_embedding_matrix) # save

```

```

1 def convert_to_ints(text, word_count, unk_count, eos=False):
2     """Convert words in text to an integer.
3     If word is not in vocab_to_int, use UNK's integer.
4     Total the number of words and UNKS.
5     Add EOS token to the end of texts"""
6     ints = []
7     for sentence in text:
8         sentence_ints = []
9         for word in sentence:
10
11             word_count += 1
12             if word in vocab_to_int:
13                 sentence_ints.append(vocab_to_int[word])
14             else:
15                 sentence_ints.append(vocab_to_int["<UNK>"])
16             unk_count += 1
17         if eos:
18             sentence_ints.append(vocab_to_int["<EOS>"])
19         ints.append(sentence_ints)
20     return ints, word_count, unk_count

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

1 # Apply convert_to_ints to clean_summaries and clean_texts
2 word_count = 0
3 unk_count = 0
4
5 int_summaries, word_count, unk_count = convert_to_ints(clean_summaries, word_count, unk_count)
6 int_texts, word_count, unk_count = convert_to_ints(clean_texts, word_count, unk_count, eos=True)

```

```

1 def create_lengths(text):
2     """Create a data frame of the sentence lengths from a text"""
3     lengths = []
4     for sentence in text:
5         lengths.append(len(sentence))
6     return pd.DataFrame(lengths, columns=['counts'])

```

```

1 lengths_summaries = create_lengths(int_summaries)
2 lengths_texts = create_lengths(int_texts)
3
4 print("Summaries:")
5 print(lengths_summaries.describe())
6 print()
7 print("Texts:")
8 print(lengths_texts.describe())

```

Summaries:

```

counts
count 29170.000000
mean   10.748166
std     2.980412
min     3.000000
25%     9.000000
50%    10.000000
75%    13.000000
max    26.000000

```

Texts:

```

counts
count 29170.000000
mean   23.781796
std     6.073380
min     5.000000
25%    19.000000
50%    23.000000
75%    28.000000
max    51.000000

```

```

1 # Inspect the length of texts
2 print(np.percentile(lengths_texts.counts, 90))
3 print(np.percentile(lengths_texts.counts, 95))
4 print(np.percentile(lengths_texts.counts, 99))

```

```

32.0
34.0
39.0

```

```

1 # Inspect the length of summaries
2 print(np.percentile(lengths_summaries.counts, 90))
3 print(np.percentile(lengths_summaries.counts, 95))
4 print(np.percentile(lengths_summaries.counts, 99))

```

```

15.0
16.0
19.0

```

```

1 def unk_counter(sentence):
2     """Counts the number of time UNK appears in a sentence."""
3     unk_count = 0
4     for word in sentence:
5         if word == vocab_to_int["<UNK>"]:
6             unk_count += 1
7     return unk_count

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

1 # Sort the summaries and texts by the length of the texts, shortest to longest
2 # Limit the length of summaries and texts based on the min and max ranges.
3 # Remove news that include too many UNKs
4
5 max_text_length = 1000
6 max_summary_length = 100
7 min_length = 2
8
9 unk_text_limit = 1
10 unk_summary_limit = 0
11
12 y = [] # summaries dataset
13 X = [] # texts dataset
14
15 for length in range(min(lengths_texts.counts), max_text_length):
16     for count, words in enumerate(int_summaries):
17         if (len(int_summaries[count]) >= min_length and
18             len(int_summaries[count]) <= max_summary_length and
19             len(int_texts[count]) >= min_length and
20             unk_counter(int_summaries[count]) <= unk_summary_limit and
21             unk_counter(int_texts[count]) <= unk_text_limit and
22             length == len(int_texts[count])):
23         ):
24             y.append(int_summaries[count])
25             X.append(int_texts[count])
26
27 # Compare lengths to ensure they match
28 print('y size: ', len(y))
29 print('X size: ', len(X))

```

```

y size: 23770
X size: 23770

```

Train-Test Split

```

train 90 %
test 10 %

```

```

1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.1)
3
4 print('Dataset size: ', len(X))
5 print('Train size:', len(X_train))
6 print('Test size:', len(X_test))

```

```

Dataset size: 23770
Train size: 21393
Test size: 2377

```

```

1 def find_max(data):
2     return max([len(d) for d in data])
3
4 def find_min(data):
5     return min([len(d) for d in data])

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

1 # data sory by X but without y sorted, maual y sort
2 print('# Dataset')
3 print("The shortest text length: ", find_min(X))
4 print("The longest text length: ", find_max(X))
5 print()
6 print("The shortest summary length: ", find_min(y))
7 print("The longest summary length : ", find_max(y))
8
9 print('-----')
10 print('# Train set')
11 print("The shortest text length: ", find_min(X_train))
12 print("The longest text length: ", find_max(X_train))
13 print()
14 print("The shortest summary length: ", find_min(y_train))
15 print("The longest summary length: ", find_max(y_train))
16
17 print('-----')
18 print('# Test set')
19 print("The shortest text length: ", find_min(X_test))
20 print("The longest text length: ", find_max(X_test))
21 print()
22 print("The shortest summary length: ", find_min(y_test))
23 print("The longest summary length: ", find_max(y_test))

```

Dataset

The shortest text length: 5
The longest text length: 49

The shortest summary length: 3
The longest summary length : 26

Train set

The shortest text length: 5
The longest text length: 49

The shortest summary length: 3
The longest summary length: 26

Test set

The shortest text length: 5
The longest text length: 45

The shortest summary length: 3
The longest summary length: 22

Save Dataset

```

1 dataset = {
2     'data': X, 'target': y,
3     'X_train': X_train, 'X_test': X_test,
4     'y_train': y_train, 'y_test': y_test
5 }
6
7 to_json('dataset.json', dataset)

```

9. โค้ดสร้างรายงาน

```

1 from matplotlib import pyplot as plt
2 import json
3
4 # config
5 batch_size = 256
6 update_check = 26
7 num_update_check_per_epoch = 4
8
9
10 def read_txt(fname):
11     with open(fname, 'r', encoding='utf-8') as file:
12         data = file.read()
13         string_data = data.split('\n')
14         try:
15             float_data = [float(s) for s in string_data]
16             return float_data
17         except ValueError:
18             return string_data
19
20
21 def read_json(fname):
22     with open(fname, 'r') as file:
23         data = file.read()
24         json_data = json.loads(data)
25         return json_data
26
27
28 def get_data(loss, limit=500):
29     return [loss[i-1] for i in range(1, len(loss)+1) if i % num_update_check_per_epoch is 0][:limit]

```

```

1 model_dict = {
2     'Adam without Dropout Model 1': 'Adam/0_1',
3     'Adam without Dropout Model 2': 'Adam/0_2',
4     'Adam Dropout 0.2': 'Adam/20',
5     'Adam Dropout 0.35': 'Adam/35',
6     'RMSProp without Dropout': 'RMSProp/0',
7     'RMSProp Dropout 0.2': 'RMSProp/20',
8     'RMSProp Dropout 0.35': 'RMSProp/35',
9 }

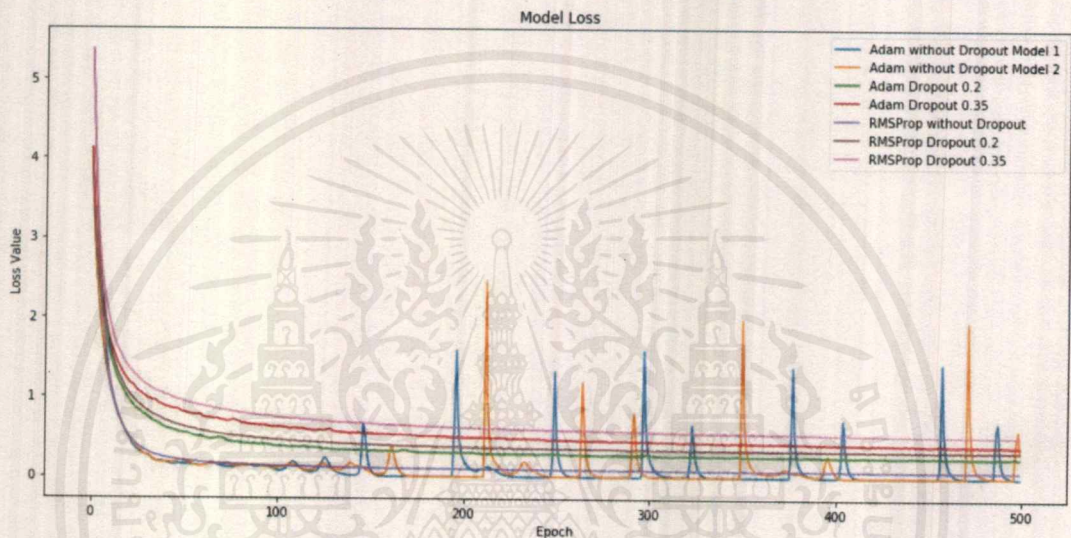
```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

1 # load loss value
2 plot_dict = {}
3 for model_name, path in model_dict.items():
4     loss = read_txt(f'{path}/summary_update_loss.txt')
5     plot_dict[model_name] = get_data(loss)
6
7
8 # plot
9 plt.figure(figsize=(15, 7))
10
11 _ = [plt.plot(loss) for loss in plot_dict.values()]
12
13 plt.title('Model Loss')
14 plt.ylabel('Loss Value')
15 plt.xlabel('Epoch')
16
17 plt.legend(plot_dict.keys())
18 plt.savefig('model_loss.png') # save to image
19 plt.show()

```



```

1 import pandas as pd
2 # load report value
3 df = pd.DataFrame()
4 for model_name, path in model_dict.items():
5     report = read_json(f'{path}/report.json')
6     report = pd.DataFrame([report])
7     df = df.append(report)
8
9 df.index = model_dict

```

```
1 df.round(4).to_csv('report.csv')
```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

10. โค้ดการนำแบบจำลองไปใช้งาน

```

1 import pythainlp
2 from pythainlp import tokenize
3 from num2words import num2words
4 import numpy as np
5 import tensorflow as tf
6 import re
7 import time
8 import json

```

```

1 print('PyThaiNLP Version: {}'.format(pythainlp.__version__))
2 print('TensorFlow Version: {}'.format(tf.__version__))

```

PyThaiNLP Version: 2.0.5
TensorFlow Version: 1.13.1

```

1 def read_json(fname, key_int=False):
2     with open(fname, 'r') as file:
3         data = file.read()
4         json_data = json.loads(data)
5
6         if not key_int:
7             return json_data
8
9         json_data = {int(key): value for key, value in json_data.items()}
10        return json_data
11

```

```

1 def read_txt(fname):
2     with open(fname, 'r', encoding='utf-8') as file:
3         data = file.read()
4         string_data = data.split('\n')
5         try:
6             float_data = [float(s) for s in string_data]
7             return float_data
8         except ValueError:
9             return string_data

```

```

1 def text_to_seq(text):
2     """Prepare the text for the model"""
3     return [vocab_to_int.get(word, vocab_to_int['<UNK>']) for word in text]

```

```

1 def contractions(text):
2     # replace contractions with their longer forms
3     for word1, word2 in contractions_dict.items():
4         text = text.replace(word1, word2)
5     return text

```

```

1 def clean_token(tokens, stopword_dict={}, use_remove_symbol=True):
2     clean = []
3     for token in tokens:
4         token = token.strip(' ') # remove white space
5         if token in stopword_dict:
6             continue
7
8         if use_remove_symbol:
9             token = remove_symbol(token)
10            if len(token) <= 1:
11                continue
12
13            clean.append(token)
14    return clean

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

1 def build_vocab(text, unk_method='replace', use_norm_num=True, **kwargs):
2     """
3     unk_method: unknow vocab method
4     replace: replace know vocab to <unk>
5     remove: pass vocab (remove vocab)
6     ignore: can not operate
7     """
8
9     # to lower-case
10    text = text.lower()
11    # replace contractions with their longer forms
12    text = contractions(text)
13    # word tokenize
14    tokens = TK.word_tokenize(text)
15
16    segments = []
17    for token in tokens:
18        if use_norm_num:
19            # convert number to word
20            token_list = norm_num(token, **kwargs)
21        else:
22            token_list = [token]
23        # append to segments list
24        segments += token_list
25    return segments

```

```

1 def norm_num(string, use_num2word=False, lang='thai'):
2     """return dtype: list"""
3
4     ignore_string = ['infinity']
5
6     if string in ignore_string:
7         return string
8
9     try:
10    new_string = string.replace('.', '')
11    num = float(new_string)
12    try:
13        num = int(new_string)
14    except ValueError:
15        pass
16
17    if use_num2word:
18        word = num2words(num, lang=lang) # to word
19        tokens = TK.word_tokenize(word)
20    else:
21        tokens = [str(num)]
22    return tokens # dtype: list
23 except ValueError:
24    return [string]

```

```

1 def remove_symbol(word):
2     # Format words and remove unwanted characters
3     word = re.sub(r'https?:\V.*[\r\n]*', "", word, flags=re.MULTILINE)
4     word = re.sub(r'\<a href', "", word)
5     word = re.sub(r'&#', "", word)
6     # word = re.sub(r'[_\-\%()|+\&=%,!\?:#\$\@|\[\]/]', "", word) # with . symbol
7     word = re.sub(r'[_\-\%()|+\&=%,!\?:#\$\@|\[\]/]', "", word) # without . symbol
8     word = re.sub(r'\<br />', "", word)
9     word = re.sub(r'\\', "", word)
10    return word

```

```

1 tk_engine = 'newmm'
2 batch_size = 256

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Load Value

```

1 vocab_to_int = read_json('vocab_to_int.json')
2 int_to_vocab = read_json('int_to_vocab.json', key_int=True)
3 word_embedding_matrix = np.load('word_embedding_matrix.npy') # load
4
5 stopwords_ = set(read_txt('stopwords.txt'))
6 DICTIONARY = set(read_txt('dictionary.txt'))
7 contractions_dict = read_json('contractions_dict.json')
8
9 TK = tokenize.Tokenizer(DICTIONARY, engine=tk_engine)

```

```
1 print('number of word output space: ', len(int_to_vocab))
```

number of word output space: 19878

```

1 def textsumm(input_sentence, checkpoint):
2     build_input_sentence = [build_vocab(text) for text in input_sentence]
3     clean_input_sentence = [clean_token(tokens, stopword_dict={}) for tokens in build_input_sentence]
4     texts = [text_to_seq(text) for text in clean_input_sentence] # text for input
5
6     loaded_graph = tf.Graph()
7     with tf.Session(graph=loaded_graph) as sess:
8         # Load saved model
9         loader = tf.train.import_meta_graph(checkpoint + '.meta')
10        loader.restore(sess, checkpoint)
11
12        input_data = loaded_graph.get_tensor_by_name('input:0')
13        logits = loaded_graph.get_tensor_by_name('predictions:0')
14        text_length = loaded_graph.get_tensor_by_name('text_length:0')
15        summary_length = loaded_graph.get_tensor_by_name('summary_length:0')
16        keep_prob = loaded_graph.get_tensor_by_name('keep_prob:0')
17
18        outputs = []
19        for text in texts:
20            #Multiply by batch_size to match the model's input parameters
21            answer_logits = sess.run(logits, {input_data: [text]*batch_size,
22                summary_length: [26],
23                text_length: [len(text)]*batch_size,
24                keep_prob: 1.0})[0]
25            outputs.append(answer_logits)
26
27        # Remove the padding
28        pad = vocab_to_int["<PAD>"]
29        sums = [[int_to_vocab[i] for i in output if i != pad] for output in outputs]
30
31        return sums

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

11. โค้ด ROUGE metric

```

1  """ROUGE metric implementation.
2  This is a modified and slightly extended version of
3  https://github.com/google/seq2seq.
4  """
5
6  import itertools
7  import numpy as np
8
9
10 def _get_ngrams(n, text):
11     """Calculates n-grams.
12     Args:
13         n: which n-grams to calculate
14         text: An array of tokens
15     Returns:
16         A set of n-grams
17     """
18     ngram_set = set()
19     text_length = len(text)
20     max_index_ngram_start = text_length - n
21     for i in range(max_index_ngram_start + 1):
22         ngram_set.add(tuple(text[i:i + n]))
23     return ngram_set
24
25
26 def _split_into_words(sentences):
27     """Splits multiple sentences into words and flattens the result"""
28     return list(itertools.chain(*[_split("|") for _ in sentences]))
29
30
31 def _get_word_ngrams(n, sentences):
32     """Calculates word n-grams for multiple sentences.
33     """
34     assert len(sentences) > 0
35     assert n > 0
36
37     words = _split_into_words(sentences)
38     return _get_ngrams(n, words)
39
40
41 def _len_lcs(x, y):
42     """
43     Returns the length of the Longest Common Subsequence between sequences x
44     and y.
45     Source: http://www.algorithmist.com/index.php/Longest\_Common\_Subsequence
46     Args:
47         x: sequence of words
48         y: sequence of words
49     Returns
50         integer: Length of LCS between x and y
51     """
52     table = _lcs(x, y)
53     n, m = len(x), len(y)
54     return table[n, m]

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

57 def _lcs(x, y):
58     """
59     Computes the length of the longest common subsequence (lcs) between two
60     strings. The implementation below uses a DP programming algorithm and runs
61     in O(nm) time where n = len(x) and m = len(y).
62     Source: http://www.algorithmist.com/index.php/Longest\_Common\_Subsequence
63     Args:
64         x: collection of words
65         y: collection of words
66     Returns:
67         Table of dictionary of coord and len lcs
68     """
69     n, m = len(x), len(y)
70     table = dict()
71     for i in range(n + 1):
72         for j in range(m + 1):
73             if i == 0 or j == 0:
74                 table[i, j] = 0
75             elif x[i - 1] == y[j - 1]:
76                 table[i, j] = table[i - 1, j - 1] + 1
77             else:
78                 table[i, j] = max(table[i - 1, j], table[i, j - 1])
79     return table

82 def _recon_lcs(x, y):
83     """
84     Returns the Longest Subsequence between x and y.
85     Source: http://www.algorithmist.com/index.php/Longest\_Common\_Subsequence
86     Args:
87         x: sequence of words
88         y: sequence of words
89     Returns:
90         sequence: LCS of x and y
91     """
92     i, j = len(x), len(y)
93     table = _lcs(x, y)
94
95     def _recon(i, j):
96         """private recon calculation"""
97         if i == 0 or j == 0:
98             return []
99         elif x[i - 1] == y[j - 1]:
100             return _recon(i - 1, j - 1) + [(x[i - 1], i)]
101         elif table[i - 1, j] > table[i, j - 1]:
102             return _recon(i - 1, j)
103         else:
104             return _recon(i, j - 1)
105
106     recon_tuple = tuple(map(lambda x: x[0], _recon(i, j)))
107     return recon_tuple

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

110 def rouge_n(evaluated_sentences, reference_sentences, n=2):
111     """
112     Computes ROUGE-N of two text collections of sentences.
113     Source: http://research.microsoft.com/en-us/um/people/cyl/download/
114     papers/rouge-working-note-v1.3.1.pdf
115     Args:
116     evaluated_sentences: The sentences that have been picked by the summarizer
117     reference_sentences: The sentences from the referene set
118     n: Size of ngram. Defaults to 2.
119     Returns:
120     A tuple (f1, precision, recall) for ROUGE-N
121     Raises:
122     ValueError: raises exception if a param has len <= 0
123     """
124     if len(evaluated_sentences) <= 0 or len(reference_sentences) <= 0:
125         raise ValueError("Collections must contain at least 1 sentence.")
126
127     evaluated_ngrams = _get_word_ngrams(n, evaluated_sentences)
128     reference_ngrams = _get_word_ngrams(n, reference_sentences)
129     reference_count = len(reference_ngrams)
130     evaluated_count = len(evaluated_ngrams)
131
132     # Gets the overlapping ngrams between evaluated and reference
133     overlapping_ngrams = evaluated_ngrams.intersection(reference_ngrams)
134     overlapping_count = len(overlapping_ngrams)
135
136     # Handle edge case. This isn't mathematically correct, but it's good enough
137     if evaluated_count == 0:
138         precision = 0.0
139     else:
140         precision = overlapping_count / evaluated_count
141
142     if reference_count == 0:
143         recall = 0.0
144     else:
145         recall = overlapping_count / reference_count
146
147     f1_score = 2.0 * ((precision * recall) / (precision + recall + 1e-8))
148
149     # return overlapping_count / reference_count
150     return f1_score, precision, recall
151
152 def f_p_r_lcs(llcs, m, n):
153     """
154     Computes the LCS-based F-measure score
155     Source: http://research.microsoft.com/en-us/um/people/cyl/download/papers/
156     rouge-working-note-v1.3.1.pdf
157     Args:
158     llcs: Length of LCS
159     m: number of words in reference summary
160     n: number of words in candidate summary
161     Returns:
162     Float. LCS-based F-measure score
163     """
164
165     r_lcs = llcs / m
166     p_lcs = llcs / n
167     beta = p_lcs / (r_lcs + 1e-12)
168     num = (1 + (beta**2)) * r_lcs * p_lcs
169     denom = r_lcs + ((beta**2) * p_lcs)
170     f_lcs = num / (denom + 1e-12)
171     return f_lcs, p_lcs, r_lcs

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

174 def rouge_l_sentence_level(evaluated_sentences, reference_sentences):
175     """
176     Computes ROUGE-L (sentence level) of two text collections of sentences.
177     http://research.microsoft.com/en-us/um/people/cyl/download/papers/
178     rouge-working-note-v1.3.1.pdf
179     Calculated according to:
180     R_lcs = LCS(X,Y)/m
181     P_lcs = LCS(X,Y)/n
182     F_lcs = ((1 + beta^2)*R_lcs*P_lcs) / (R_lcs + (beta^2) * P_lcs)
183     where:
184     X = reference summary
185     Y = Candidate summary
186     m = length of reference summary
187     n = length of candidate summary
188     Args:
189     evaluated_sentences: The sentences that have been picked by the summarizer
190     reference_sentences: The sentences from the referene set
191     Returns:
192     A float: F_lcs
193     Raises:
194     ValueError: raises exception if a param has len <= 0
195     """
196     if len(evaluated_sentences) <= 0 or len(reference_sentences) <= 0:
197         raise ValueError("Collections must contain at least 1 sentence.")
198     reference_words = _split_into_words(reference_sentences)
199     evaluated_words = _split_into_words(evaluated_sentences)
200     m = len(reference_words)
201     n = len(evaluated_words)
202     lcs = _len_lcs(evaluated_words, reference_words)
203     return _f_p_r_lcs(lcs, m, n)
206 def union_lcs(evaluated_sentences, reference_sentence):
207     """
208     Returns LCS_u(r_i, C) which is the LCS score of the union longest common
209     subsequence between reference sentence r_i and candidate summary C. For example
210     if r_i = w1 w2 w3 w4 w5, and C contains two sentences: c1 = w1 w2 w6 w7 w8 and
211     c2 = w1 w3 w8 w9 w5, then the longest common subsequence of r_i and c1 is
212     "w1 w2" and the longest common subsequence of r_i and c2 is "w1 w3 w5". The
213     union longest common subsequence of r_i, c1, and c2 is "w1 w2 w3 w5" and
214     LCS_u(r_i, C) = 4/5.
215     Args:
216     evaluated_sentences: The sentences that have been picked by the summarizer
217     reference_sentence: One of the sentences in the reference summaries
218     Returns:
219     float: LCS_u(r_i, C)
220     ValueError:
221     Raises exception if a param has len <= 0
222     """
223     if len(evaluated_sentences) <= 0:
224         raise ValueError("Collections must contain at least 1 sentence.")
225
226     lcs_union = set()
227     reference_words = _split_into_words([reference_sentence])
228     combined_lcs_length = 0
229     for eval_s in evaluated_sentences:
230         evaluated_words = _split_into_words([eval_s])
231         lcs = set(_recon_lcs(reference_words, evaluated_words))
232         combined_lcs_length += len(lcs)
233         lcs_union = lcs_union.union(lcs)
234
235     union_lcs_count = len(lcs_union)
236     union_lcs_value = union_lcs_count / combined_lcs_length
237     return union_lcs_value

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

240 def rouge_l_summary_level(evaluated_sentences, reference_sentences):
241     """
242     Computes ROUGE-L (summary level) of two text collections of sentences.
243     http://research.microsoft.com/en-us/um/people/cyl/download/papers/
244     rouge-working-note-v1.3.1.pdf
245     Calculated according to:
246      $R_{lcs} = \text{SUM}(1, u)[\text{LCS}(\cup(r_i, c))]/m$ 
247      $P_{lcs} = \text{SUM}(1, u)[\text{LCS}(\cup(r_i, c))]/n$ 
248      $F_{lcs} = ((1 + \text{beta}^2) * R_{lcs} * P_{lcs}) / (R_{lcs} + (\text{beta}^2) * P_{lcs})$ 
249     where:
250     SUM(i,u) = SUM from i through u
251     u = number of sentences in reference summary
252     C = Candidate summary made up of v sentences
253     m = number of words in reference summary
254     n = number of words in candidate summary
255     Args:
256     evaluated_sentences: The sentences that have been picked by the summarizer
257     reference_sentence: One of the sentences in the reference summaries
258     Returns:
259     A float: F_lcs
260     Raises:
261     ValueError: raises exception if a param has len <= 0
262     """
263     if len(evaluated_sentences) <= 0 or len(reference_sentences) <= 0:
264         raise ValueError("Collections must contain at least 1 sentence.")
265
266     # total number of words in reference sentences
267     m = len(_split_into_words(reference_sentences))
268
269     # total number of words in evaluated sentences
270     n = len(_split_into_words(evaluated_sentences))
271
272     union_lcs_sum_across_all_references = 0
273     for ref_s in reference_sentences:
274         union_lcs_sum_across_all_references += _union_lcs(evaluated_sentences,
275                                                         ref_s)
276     return f_p_r_lcs(union_lcs_sum_across_all_references, m, n)

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

280 def rouge(hypotheses, references):
281     """Calculates average rouge scores for a list of hypotheses and
282     references"""
283
284     # Filter out hyps that are of 0 length
285     # hyps_and_refs = zip(hypotheses, references)
286     # hyps_and_refs = [_ for _ in hyps_and_refs if len(_[0]) > 0]
287     # hypotheses, references = zip(*hyps_and_refs)
288
289     # Calculate ROUGE-1 F1, precision, recall scores
290     rouge_1 = [
291         rouge_n([hyp], [ref], 1) for hyp, ref in zip(hypotheses, references)
292     ]
293     rouge_1_f, rouge_1_p, rouge_1_r = map(np.mean, zip(*rouge_1))
294
295     # Calculate ROUGE-2 F1, precision, recall scores
296     rouge_2 = [
297         rouge_n([hyp], [ref], 2) for hyp, ref in zip(hypotheses, references)
298     ]
299     rouge_2_f, rouge_2_p, rouge_2_r = map(np.mean, zip(*rouge_2))
300
301     # Calculate ROUGE-L F1, precision, recall scores
302     rouge_l = [
303         rouge_l_sentence_level([hyp], [ref])
304         for hyp, ref in zip(hypotheses, references)
305     ]
306     rouge_l_f, rouge_l_p, rouge_l_r = map(np.mean, zip(*rouge_l))
307
308     return {
309         "rouge_1/f_score": rouge_1_f,
310         "rouge_1/r_score": rouge_1_r,
311         "rouge_1/p_score": rouge_1_p,
312         "rouge_2/f_score": rouge_2_f,
313         "rouge_2/r_score": rouge_2_r,
314         "rouge_2/p_score": rouge_2_p,
315         "rouge_l/f_score": rouge_l_f,
316         "rouge_l/r_score": rouge_l_r,
317         "rouge_l/p_score": rouge_l_p,
318     }

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหาและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้