

## preprocess ข้อมูลดิบ.

ข้อมูลดิบ

อ่านไฟล์

sentiment	review
~~~~~	positive
	negative

เปลี่ยน label positive → 1  
negative → 0

sentiment	review
~~~~~	1
	0

lower case ตัวอักษร +  
ตัด tag html

ตัดคำด้วยเรื่องว่าง + ลบ stopwords ใน eng

lemmatization . ลดรูปคำเป็น base ของคำๆนั้น.

เอาแต่ละคำมา join กลับเป็นประโยค.

save ลงไฟล์ lemma\_result.csv.

cleaned_review	Label
~~~~~	1
	0

สร้าง model.

อ่านไฟล์ที่ผ่านการ preprocess

( lemma\_result.csv )

—————> Tokenizer.

โดยกำหนดจำนวนคำที่สนใจ = max-features.



fit\_on\_text.

แปลงข้อมูล "ปลายๆประโยค" จากใน column 0  
ให้เป็นคำๆ แยกกับตัววรรคและกำหนด index ให้แต่ละ  
คำโดย "เรียงตามความถี่"



text\_to\_sequences.

แปลงประโยคให้เป็น Sequence. จากนั้น. หารความยาวของ  
ประโยคและหาค่าสูงสุด. (maxlen) เนื่องจากนั้น  
เก็บ padding ให้ทุกประโยคมีความยาวเท่ากัน โดยเติม 0.  
ด้านหน้าหรือด้านหลังประโยค.



นำข้อมูลมาแยกเป็นฝึก train และ test

3 : 1.

กำหนด random\_state เป็นค่าคงที่เพื่อให้ random ไม่สุ่มเดิม.



## สร้าง model โดยกำหนดให้

### - Embedding

- มี max-feature เป็น 7000
- กำหนด dimension = 150
- กำหนดจำนวน input node = ขนาดความยาวของ sentiment.

### - Dropout แรก เป็น 0.2

### - LSTM layer โดย

- dropout ภายในเป็น 0.2
- recurrent dropout กรณีซ่อนด้วยตัว output ก็คือ  
= 0.2
- hidden node เป็น 200

### - Dropout ครั้งที่ 2 เป็น 0.2

### - Dense unit เป็น 1 มา class output ใช้ activation function เป็น 'sigmoid'