

# Information Extraction

# Data ข้อมูล —> Information ความรู้



22 ก.ค. 2018

★★★★★ Quality Review

## ร้านในตำนานอักษรศาสตร์ จุฬาฯ

ขอยกให้ร้านนี้เป็นร้านในตำนานอักษรศาสตร์ จุฬาฯ ^^ เราทานมาตั้งแต่สมัยเรียนปริญญาตรี เหมือนนานมากมาแล้ว และห่างหายไปนาน จนล่าสุดมาเจอในแอป NOW Food Delivery โดยบังเอิญ เลยลองสั่งมาทานค่ะ

เซตที่สั่งมาเป็นข้าวเหนียวไก่ทอด+ข้าวเหนียว ราคาตามในแอปคือ 33 บาท ใส่ในกล่องกระดาษวัสดุอย่างดี แพ็กมาสวยงามน่าทานตามรูปเลยคะ โดยรสชาติความอร่อยยังเหมือนเดิม ทำให้หายคิดถึงไปได้มากเลยคะ



### เมนูที่แนะนำโดยสมาชิก

ข้าวเหนียวไก่ทอด 9

ข้าวเหนียวไก่ทอด+เอ็นไก่ทอด 2

ข้าวเหนียวหมู-ไก่ 1

ส้มตำไข่เค็ม 1

ยำไก่ทอด 1

ข้าวเหนียวเนื้อ 1

ดูทั้งหมด >>

# ภาษาจัดเป็นข้อมูลแบบไม่มีโครงสร้าง

- เปิดเพลงอะไรก็ได้ของ  
ปาล์มมีต่อนหกโมงเย็น

↓                      ↓

Artist ID	Artist Name
1	Atom ชนกันต์
<u>2</u>	<u>Palmy</u>
3	Stamp
4	แจนจิ้ง



18:00 น.

## Unstructured Data



Rockbox Brick เป็นลำโพงไร้สายที่มีเบสขนาดใหญ่ที่มีรูปร่างอิฐแบบคลาสสิก

ใช้บลูทูธเพื่อเชื่อมต่อแบบไร้สายกับอุปกรณ์

สามารถเชื่อมต่อเข้ากับ โทรศัพท์ แท็บเล็ตหรือ notebook

สามารถเป็น Powerbank แบตเตอรี่ที่มีกำลังไฟ 4000 mAh

ฟังเพลงต่อเนื่อง 20 ชั่วโมงจากการชาร์จไฟครั้งเดียวไฟ

ขนาด 15.5 x 5.9 x 5.9 ซม.

ในชุดประกอบด้วย

Rockbox Brick

Micro-USB charging cable

3.5 mm audio cable พบสินค้าเพิ่มเติมจาก FRESHN REBEL

- สี : Indigo
- กำลังไฟฟ้า (วัตต์) : 4000 mAh

## Structured Data

Brand = FRESHN REBEL

Color = Indigo

Type = Portable Speaker

Bluetooth Speaker

Home Decor

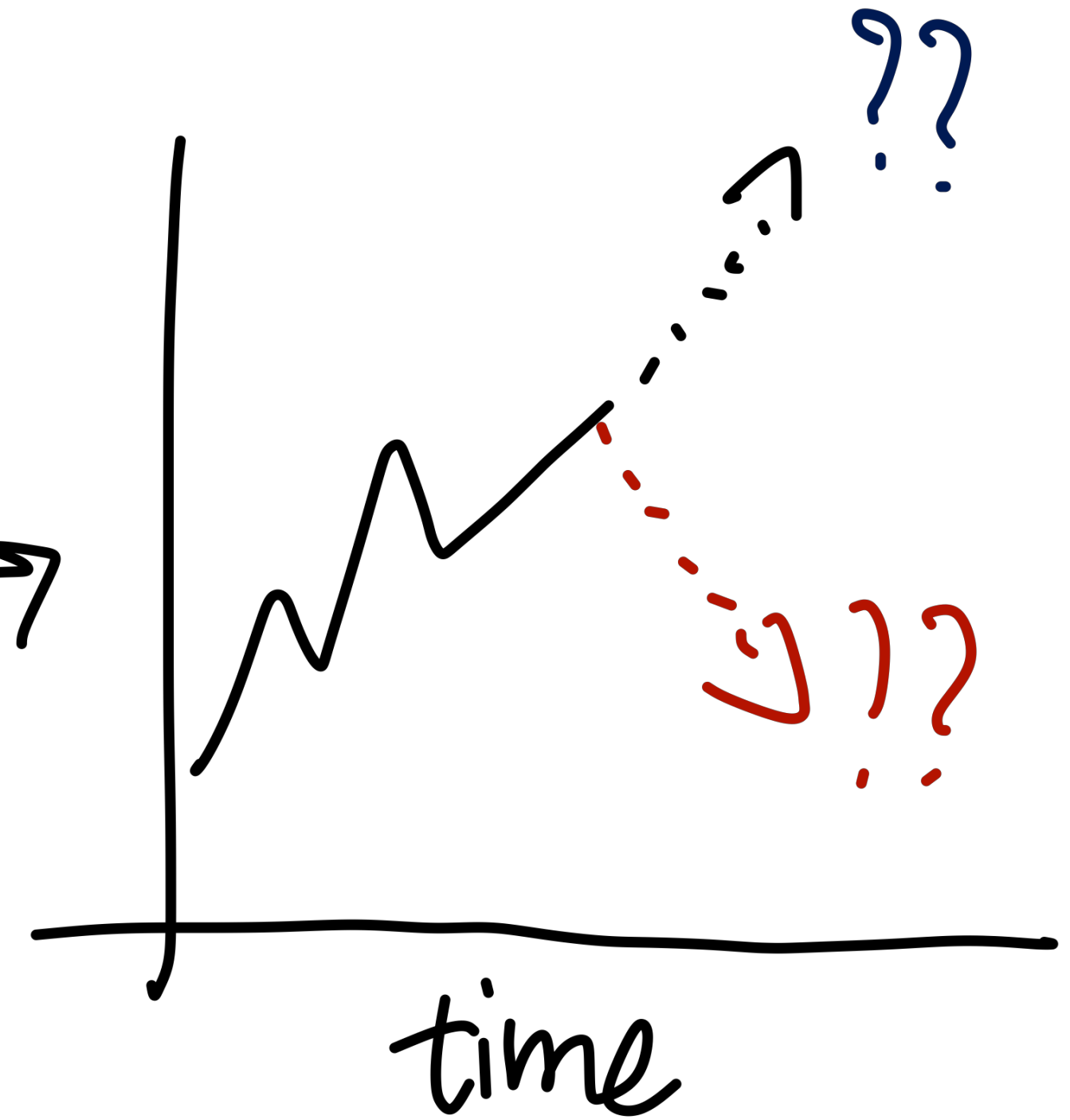
# Finances



financial report



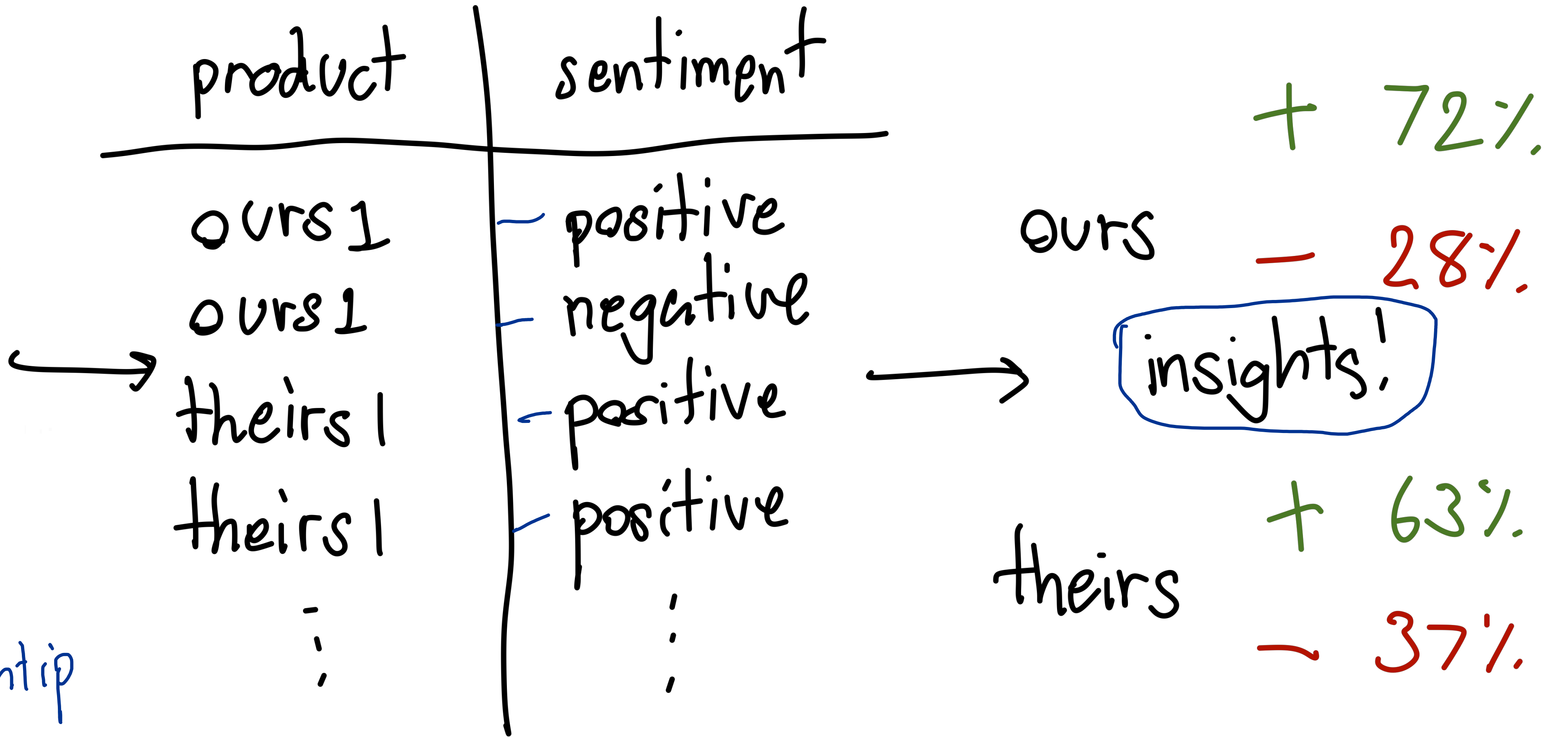
Company	event
APL	increase
UTD	unite
SET	fall
⋮	
⋮	



# Marketing + Brand Monitoring



web  
review  
facebook  
partip



# Election Forecast



newspaper  
partip  
facebook



where	candidate	sentiment
BKK	8	neutral
BKK	8	positive
BKK	6	positive
CNX	7	neutral
CNX	1	negative
⋮	⋮	



campaign  
at CNX

# Drug Administration



Clinical notes



patient	time	event	diseases
1	0	x-ray ✓	-
1	1	drug1 ✓	pneumonia ✓
2	0	MRI ✓	-
2	1	x-ray ✓	-
2	2	drug2 ✓	infection ✓



# สกัดอะไรได้บ้าง

- ชื่อคน สถานที่ทางภูมิศาสตร์ ร้านค้า ชื่อองค์กร
- ชื่อเพลง ชื่อศิลปิน ชื่ออัลบั้ม } IoT
- วัน เวลา วันที่ เหตุการณ์
- ยีนส์ โปรตีน ชื่อยา อาการทางแพทย์ เครื่องมือการวินิจฉัย ชื่อเชื้อโรค  
ชื่อโรค

# การสกัดความรู้ (Information Extraction)

- การเปลี่ยน unstructured data (data ที่เป็น text นำไปใช้ได้ยาก) เป็น structured data (data ที่เป็นตารางสามารถนำไปใช้ได้ง่าย)

# Part-of-Speech Tagging

PN

V

Noun

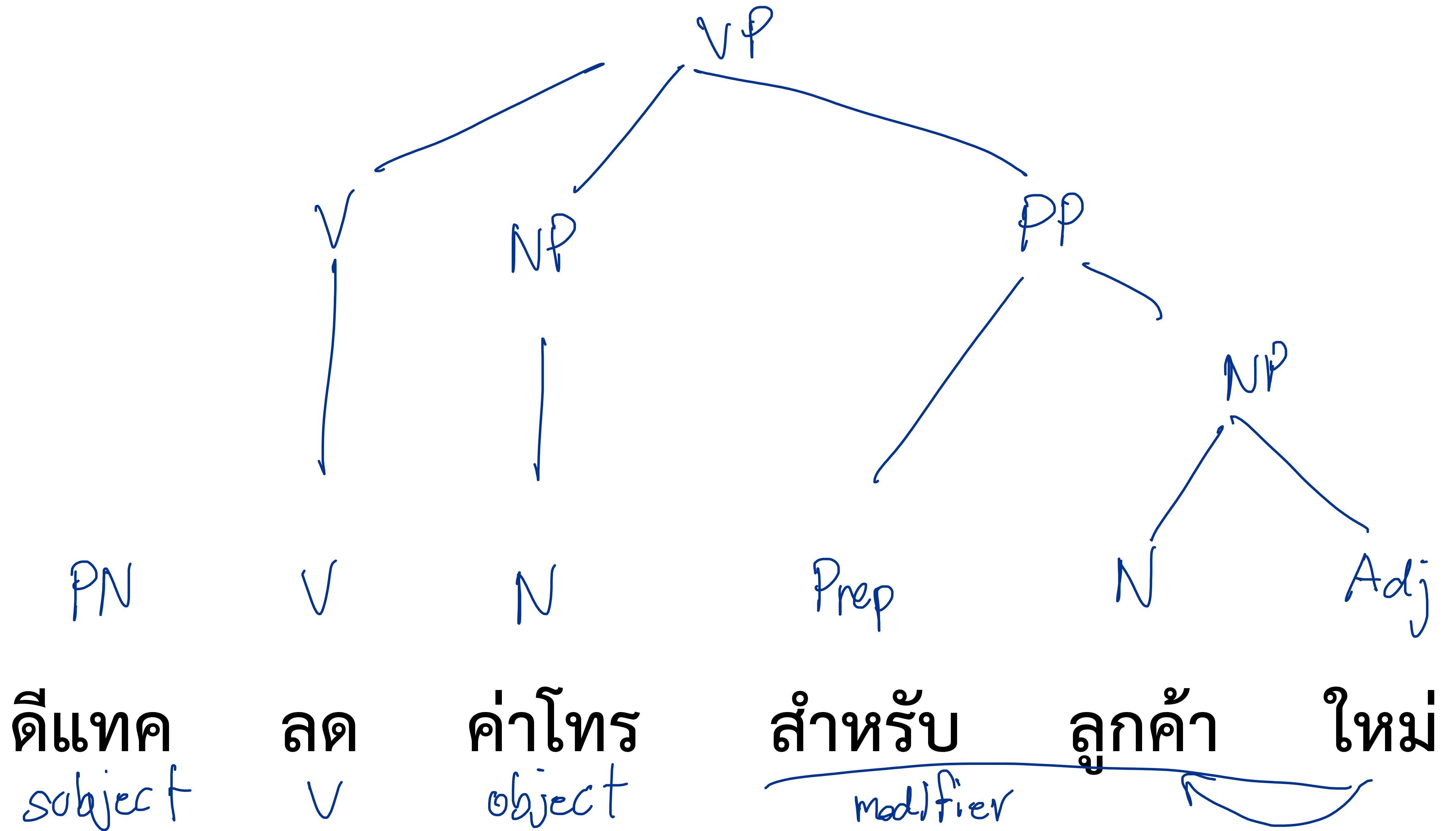
ดีแท้ค

ลค

คำโทร

Subject

Object



# Universal POS Tag

Open-class words

- ADJ
- ADV
- INTJ
- NOUN
- PROPN
- VERB

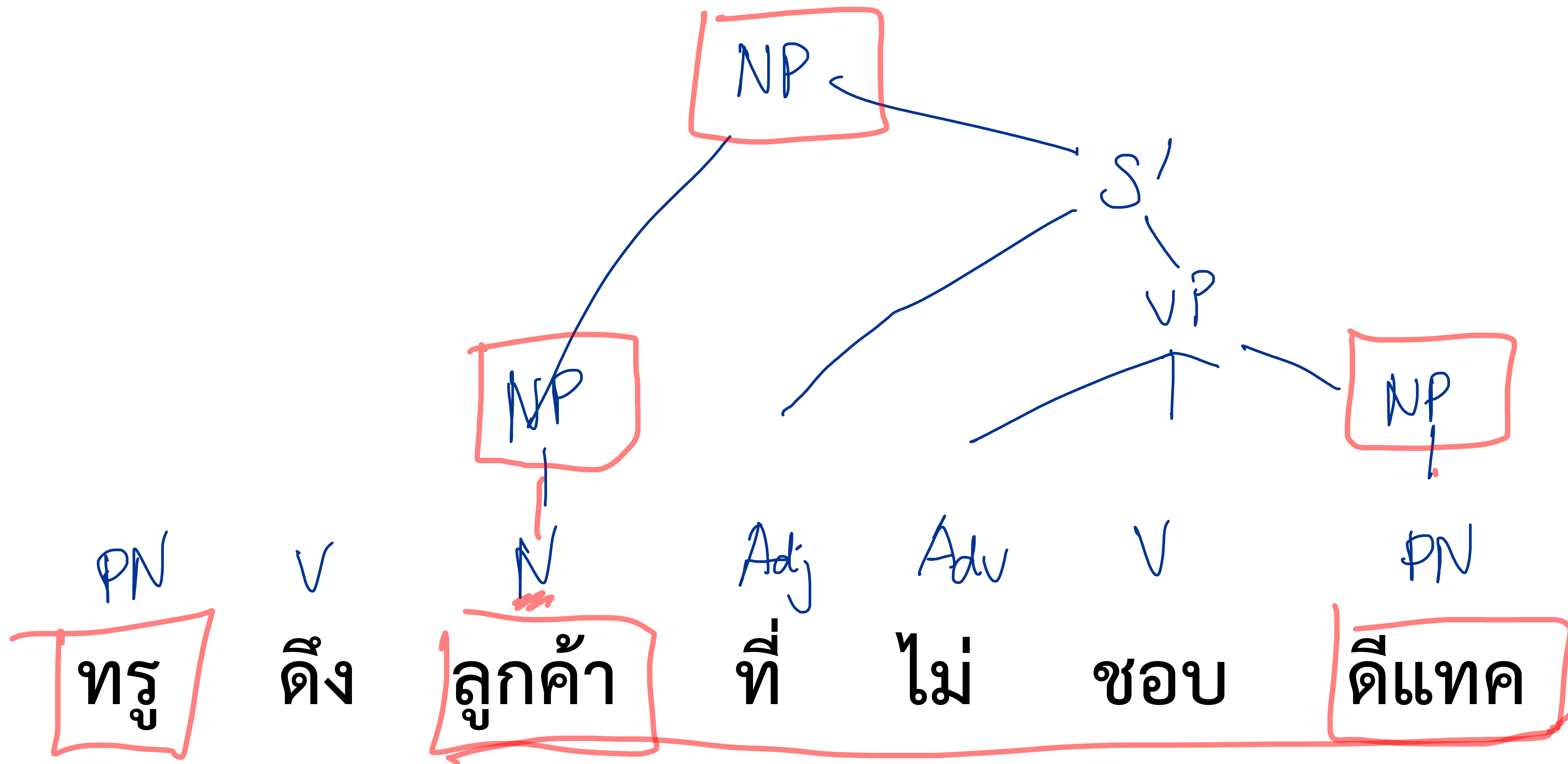
*Content word*

Closed-class words

- ADP — *preposition*
- AUX
- CCONJ
- DET — *a, an, the, my, which*
- NUM ~
- PART —
- PRON — *he, she, they*
- SCONJ

*function word*

# Base NP chunk



1) PN \*

2) N Adj

3) N

PN  
-

V  
-

N  
-

Adj  
ใหม่

Adj

Adv

V

PN

ทรู  
-

ตั้ง  
-  
V

ลูกค้า  
-

ที่  
-

ไม่

ชอบ

ดีแท้ค

Sequence Labeling



# Part-of-Speech Tagging + Base NP

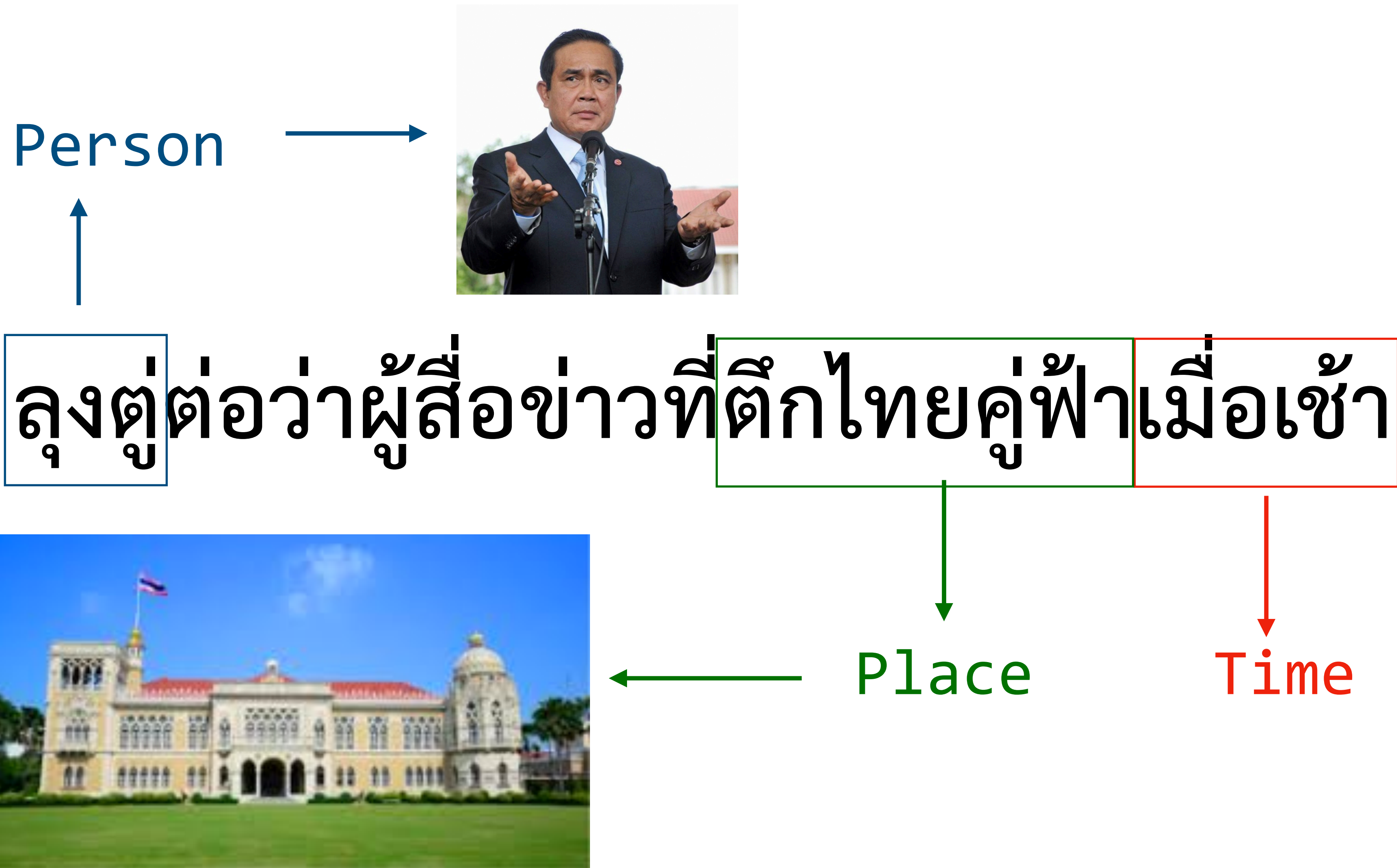
- Sequence labeling task
- การแปะชนิดของคำทำให้เราเข้าถึงความหมายได้ระดับหนึ่ง
- Base NP Chunking ช่วยสกัดความรู้เกี่ยวกับ คน สัตว์ สิ่งของ สถานที่ และสิ่งนามธรรมอื่นๆ

# Named-Entity Recognition (NER)

# displaCy Named Entity Visualizer

When **Sebastian Thrun** **PERSON** started working on self-driving cars at **Google** **ORG** in **2007** **DATE**, few people outside of the company took him seriously. “I can tell you very senior CEOs of major **American** **NORP** car companies would shake my hand and turn away because I wasn’t worth talking to,” said **Thrun** **PERSON**, now the co-founder and CEO of online higher education startup Udacity, in an interview with **Recode** **ORG** **earlier this week** **DATE**.

# Named-Entity Recognition

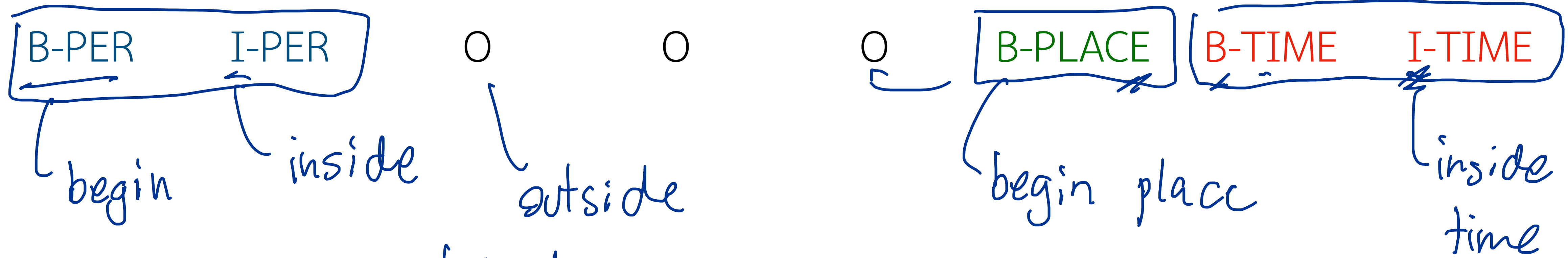


# IOB tagging NER Task Formulation BI\*

ลุงตุ๋นต่อว่าผู้สื่อข่าวที่ตึกไทยคู่ฟ้าเมื่อเช้า

O I-Time I-PER

ลุง	ตุ๋น	ต่อว่า	ผู้สื่อข่าว	ที่	ตึกไทยคู่ฟ้า	เมื่อ	เช้า
Noun	PNoun	Verb	Noun	Adj	PNoun	ADP	Noun



sequence labeling

# ชื่อคนไทย

- สนธยา คุณปลื้ม

- ก๊อบบวัน เส้น

- ปรัชญา เศรษฐกิจพอเพียง

- วัฒนา กระโปรงแดง



# ชื่อฝรั่ง

- Are you Rich?

- Play a song by Ke\$ha

- A new mission to Mars

- Dean's new single



# ชื่อยีนส์ ชื่อโปรตีน

- rpsL = ribosomal protein, small S12
- polA = DNA polymerase I
- gal = galactose
- cat = chloramphenicol resistance
- amp
- azi



# ชื่อยา

- quertiapine = Seroquel XR
- PN = penicillin != pneumonia
- IUPAC = 7-{4-[4-(2,3-dichlorophenyl) piperazin-1-yl]  
butoxy}-3,4-dihydroquinolin-2(1H)-one
- loop, potassium-sparing and thiazide diuretics  
(Dai et al, 2017)

# การรู้จำเอ็นทีตี

- NER มักถูกแก้ด้วย sequence labeling model โดยใช้ IOB label
- ยังจำเป็นต้งนึกถึงธรรมชาติของข้อมูลว่าประหลาดอย่างไร

# Sequence Labeling Model

# Sequence Labeling vs Classification

ลุง	ตู้	<u>ต่อว่า</u>	ผู้สื่อข่าว	ที่	ตึกไทยคู่ฟ้า	เมื่อ	เช้า
Noun	<u>PNoun</u>	Verb	Noun	Adj	PNoun	ADP	Noun





# Sequence ของหน่วยทางภาษาต่าง ๆ

เรื่องนี้คนแสดงนำหล่อ แต่เนื้อเรื่องน่าเบื่อ จบได้จืดมาก

positive

-

negative

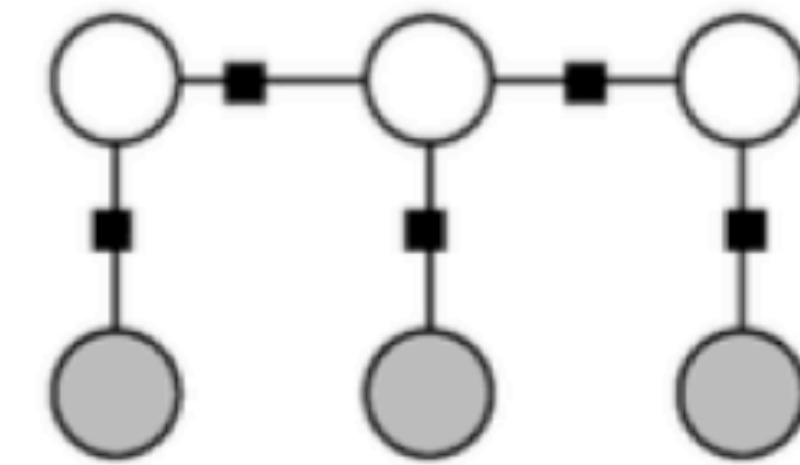
negative

# Sequence Model ที่ฮิตอยู่ขณะนี้

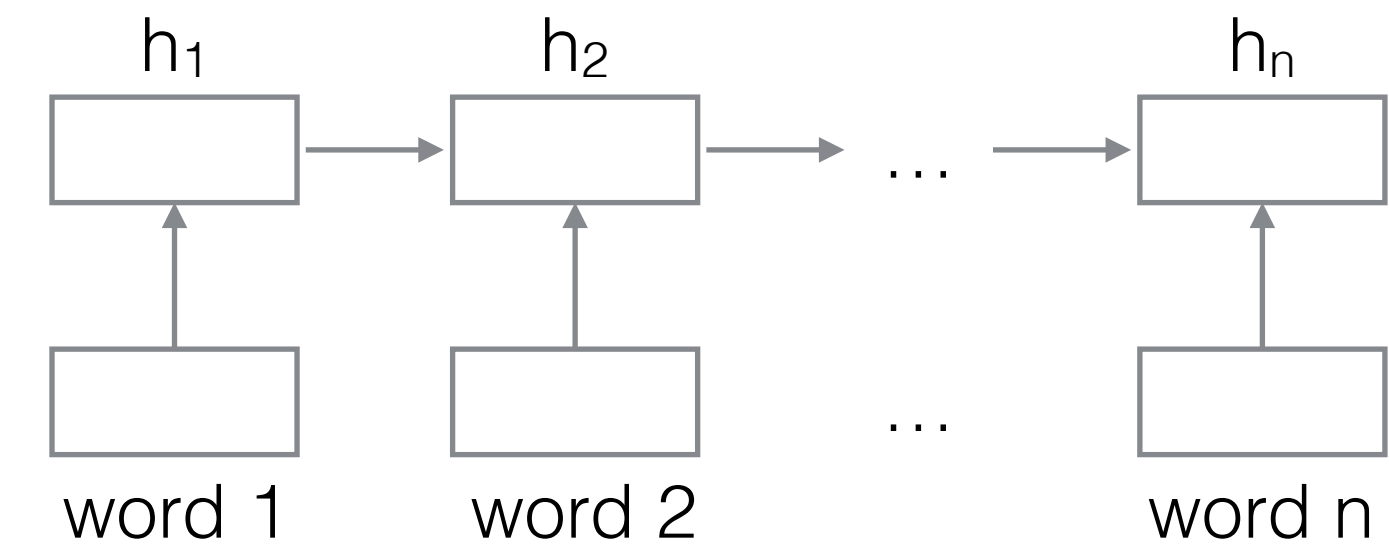
MaxEnt

Linear-chain

- Conditional Random Fields (CRF)



- Recurrent Neural Network (RNN)



- RNN + CRF

Feedforward



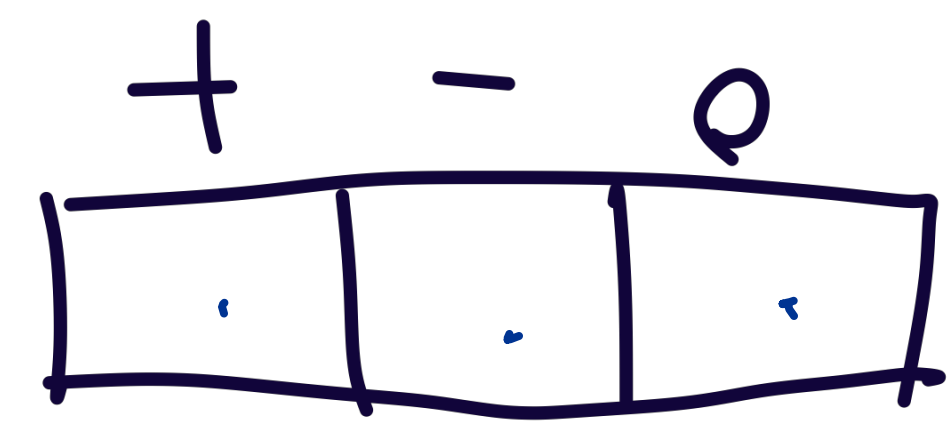
# Sequence Labeling Model

- ภาษาที่มีการเรียงตัวกันเป็นลำดับ
- ผลจะแม่นยำขึ้นถ้า Label มีความเกี่ยวเนื่องกันใน sequence  
และจำนวน Label = จำนวนหน่วย

# Conditional Random Fields (CRF)

# Logistic Regression

MaxEnt

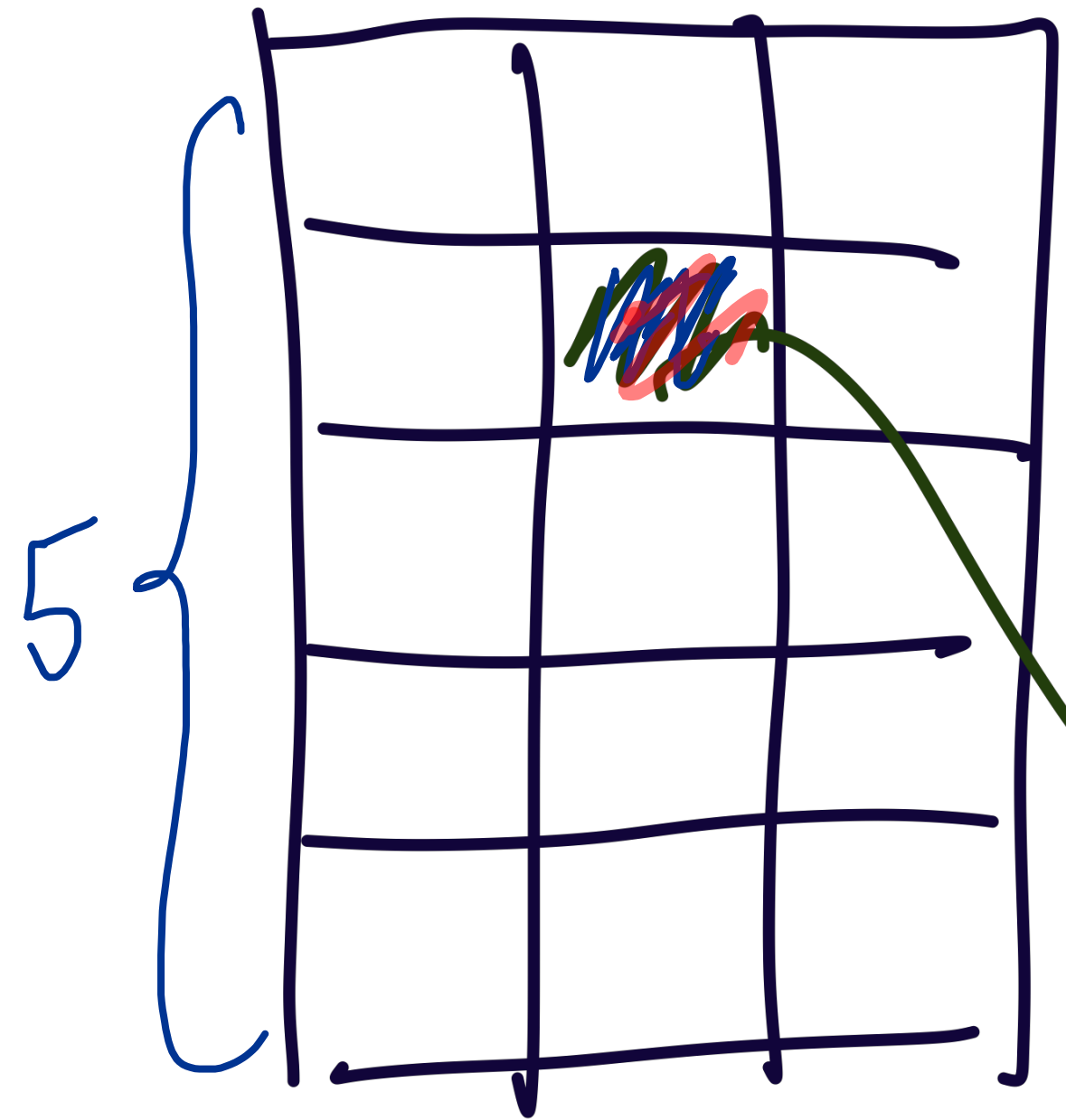


bias

+

-

0



parameter matrix

feature  $w_{ij}$   
กับ -  $w_{ij}$

5 x 3

CRF

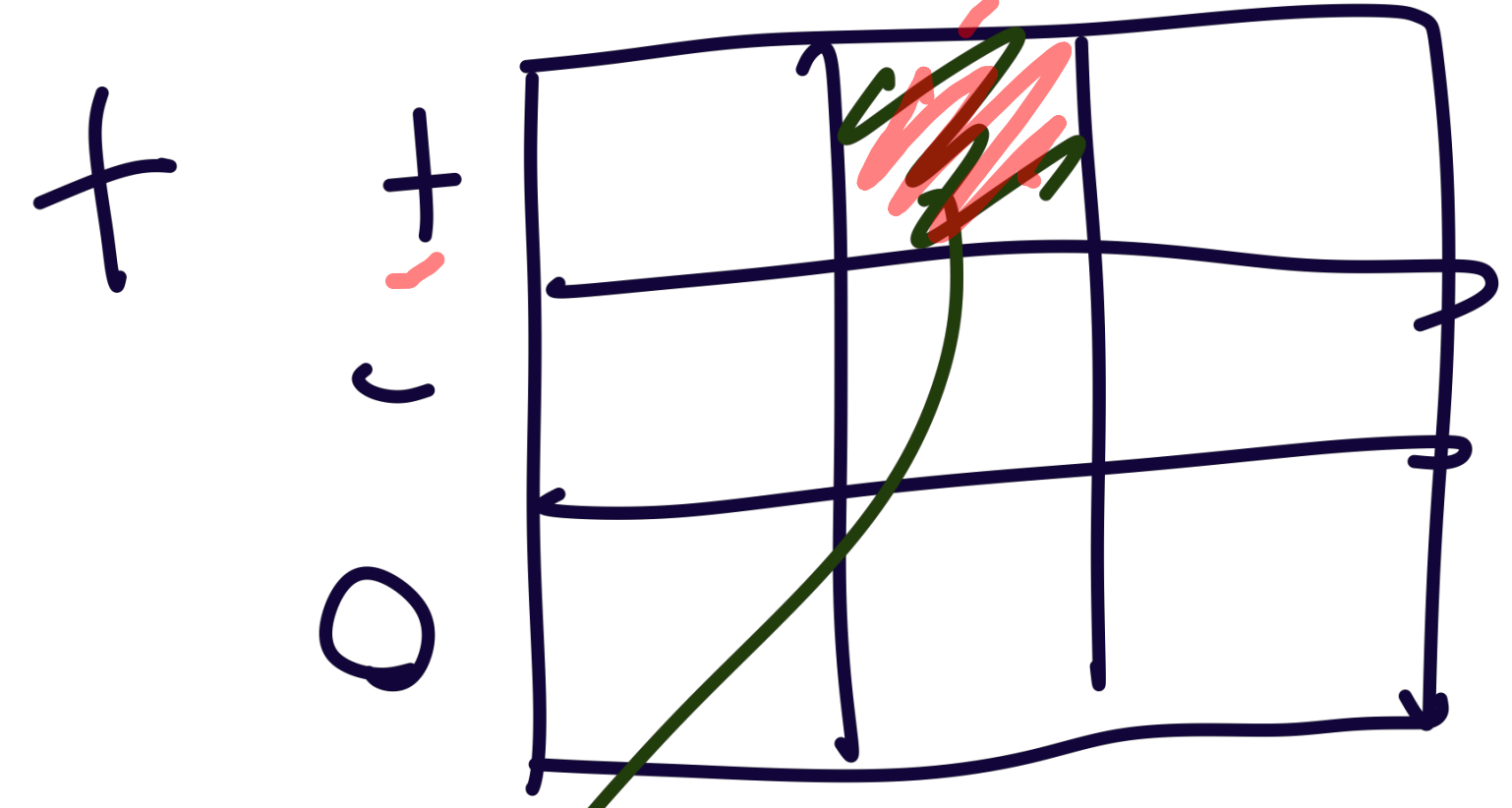
transition

label compatibility

+

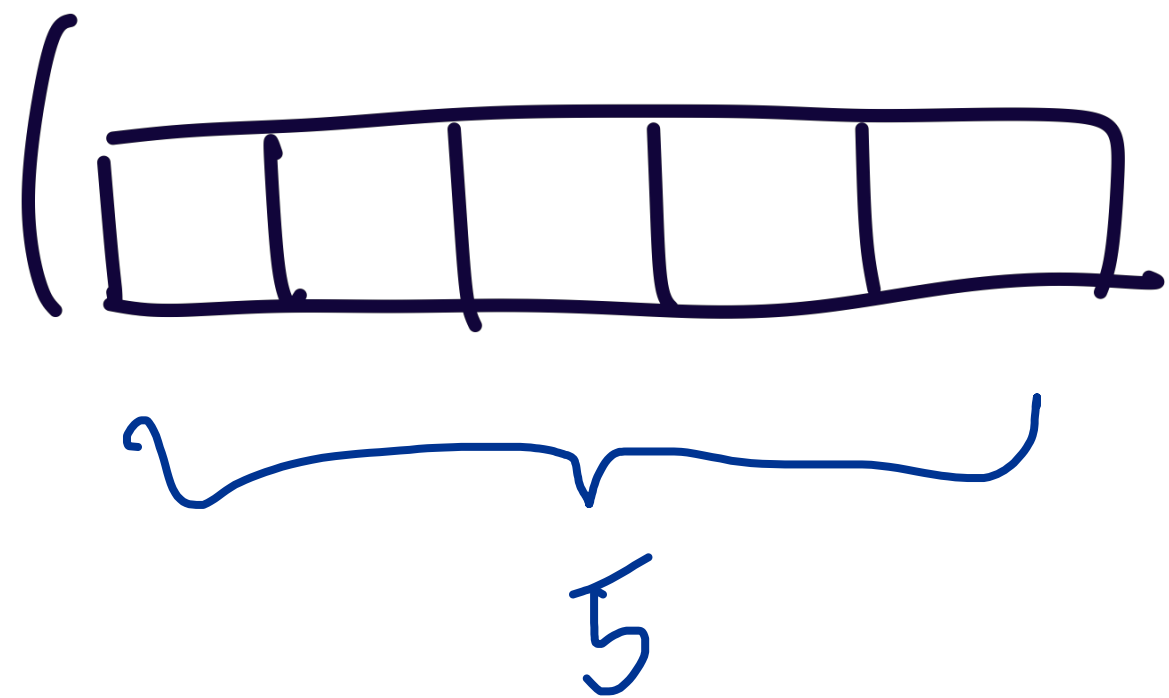
-

0

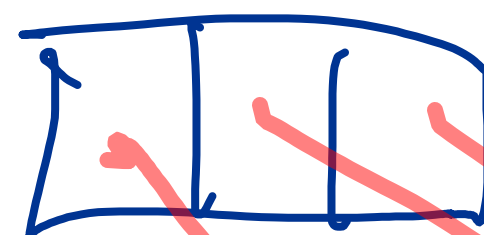
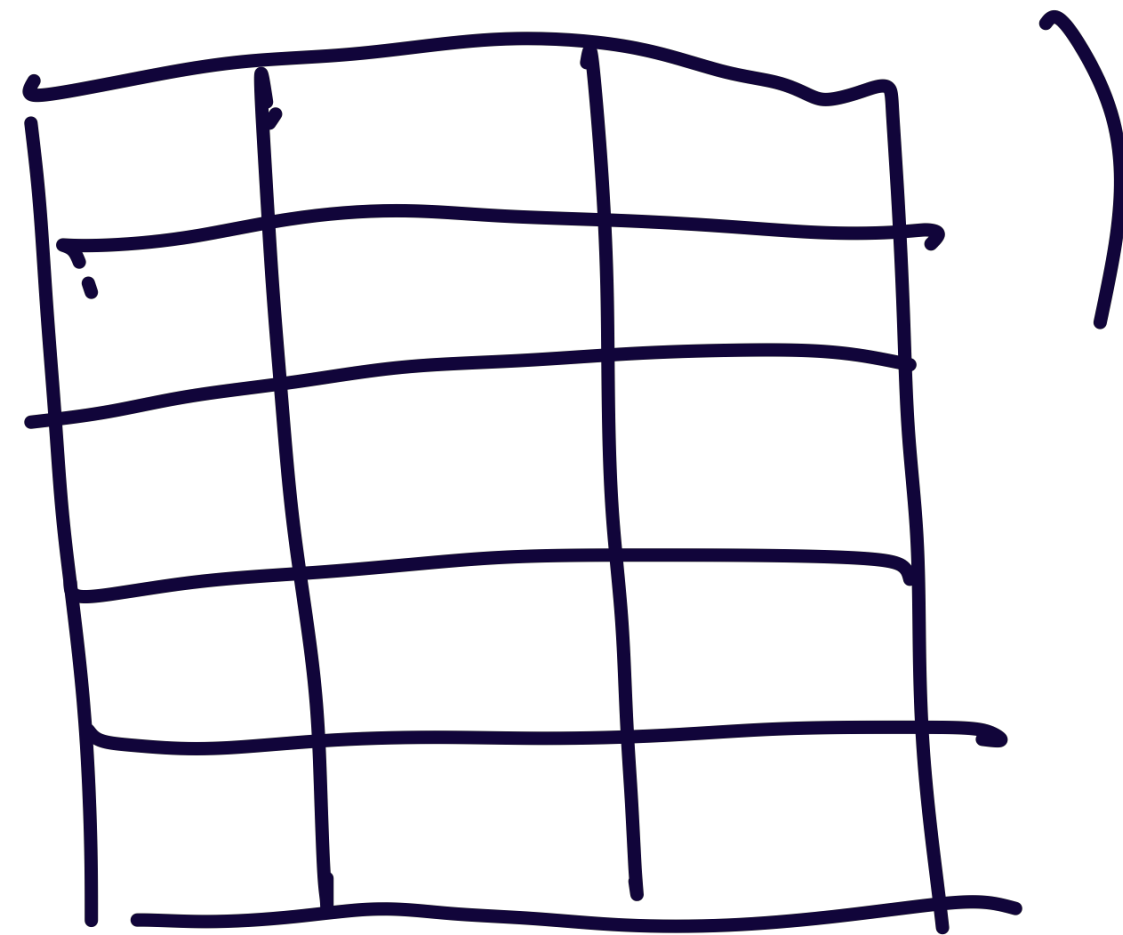


+ สถานด้วย - หนีด้วย  
ค่าบวก +, -

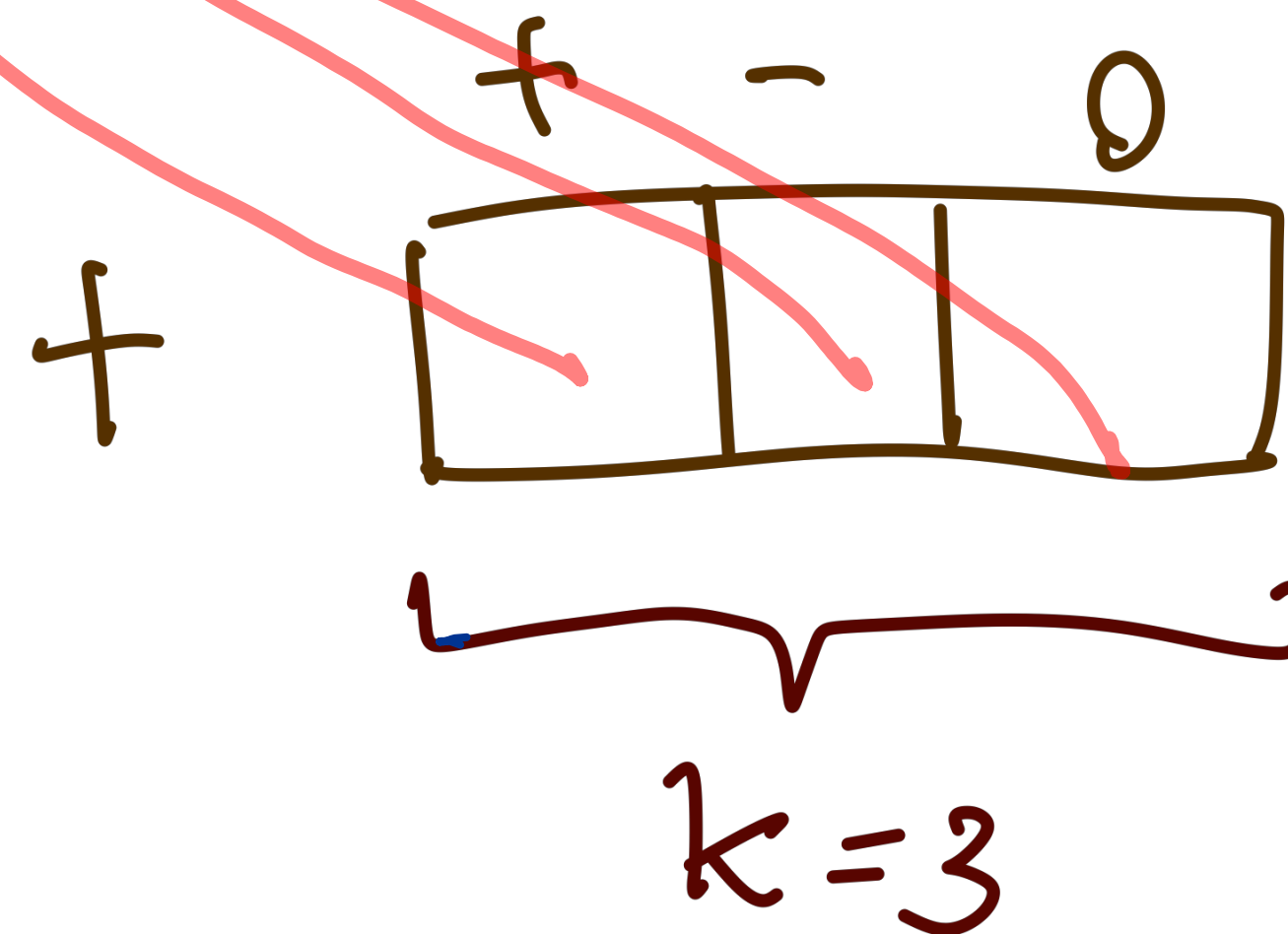
feature vector



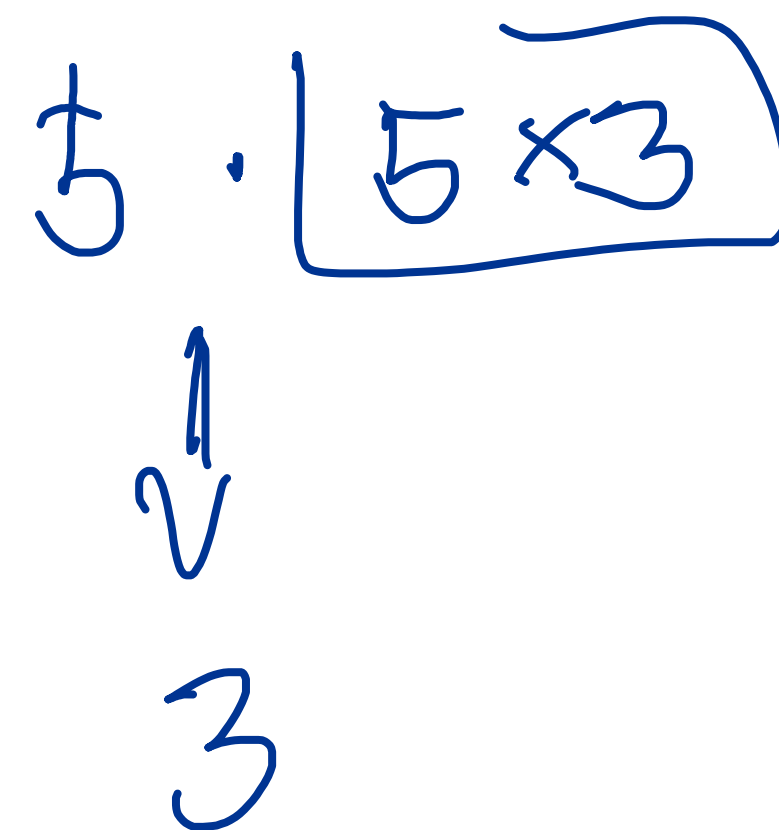
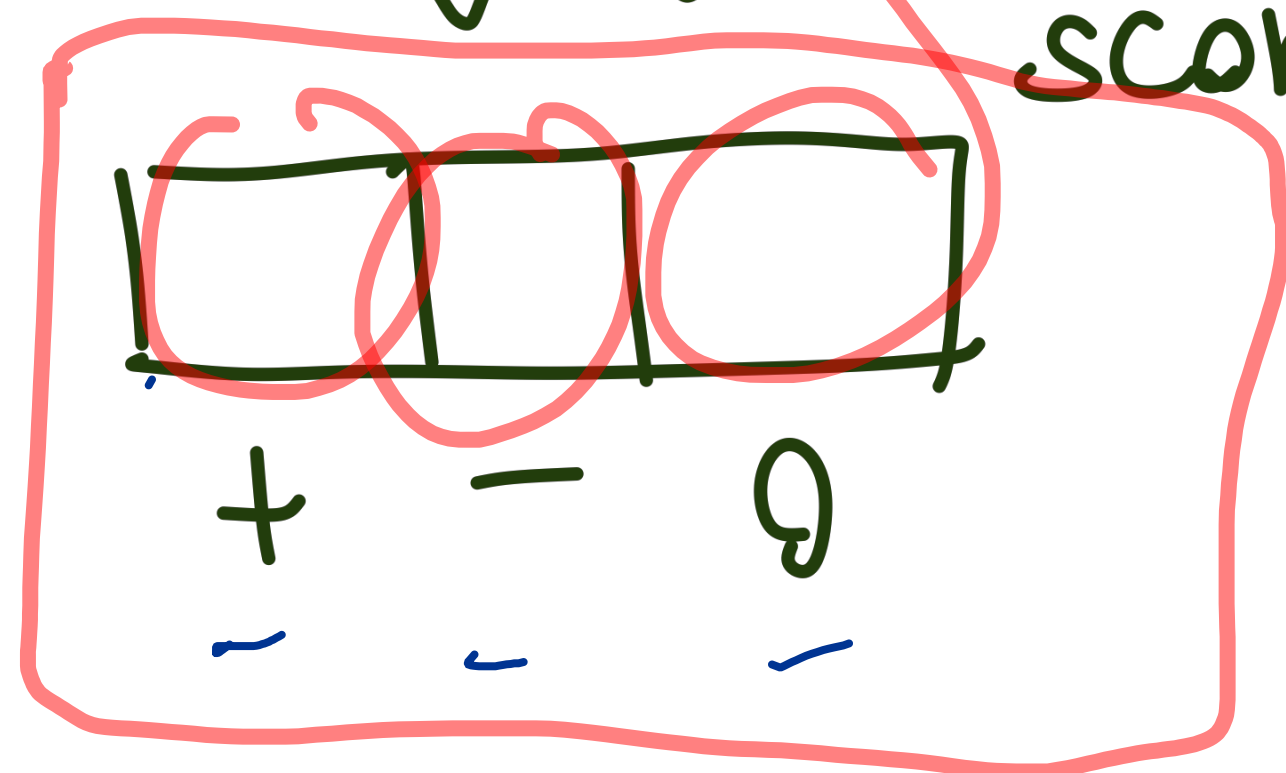
parameter matrix



bias



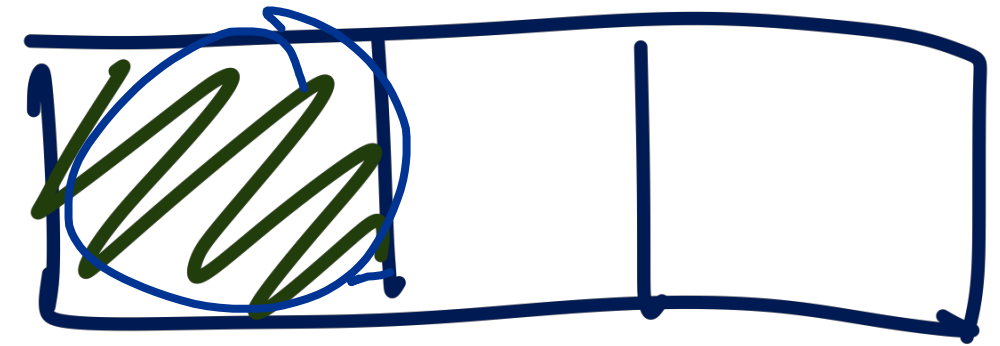
unnormalized score



t=0

sentence 1

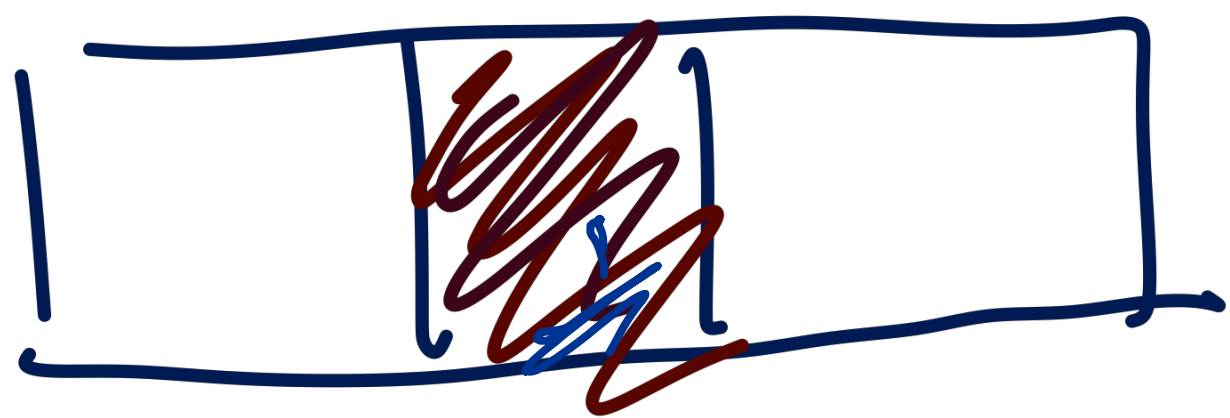
+ - 0



t=1

sentence 2

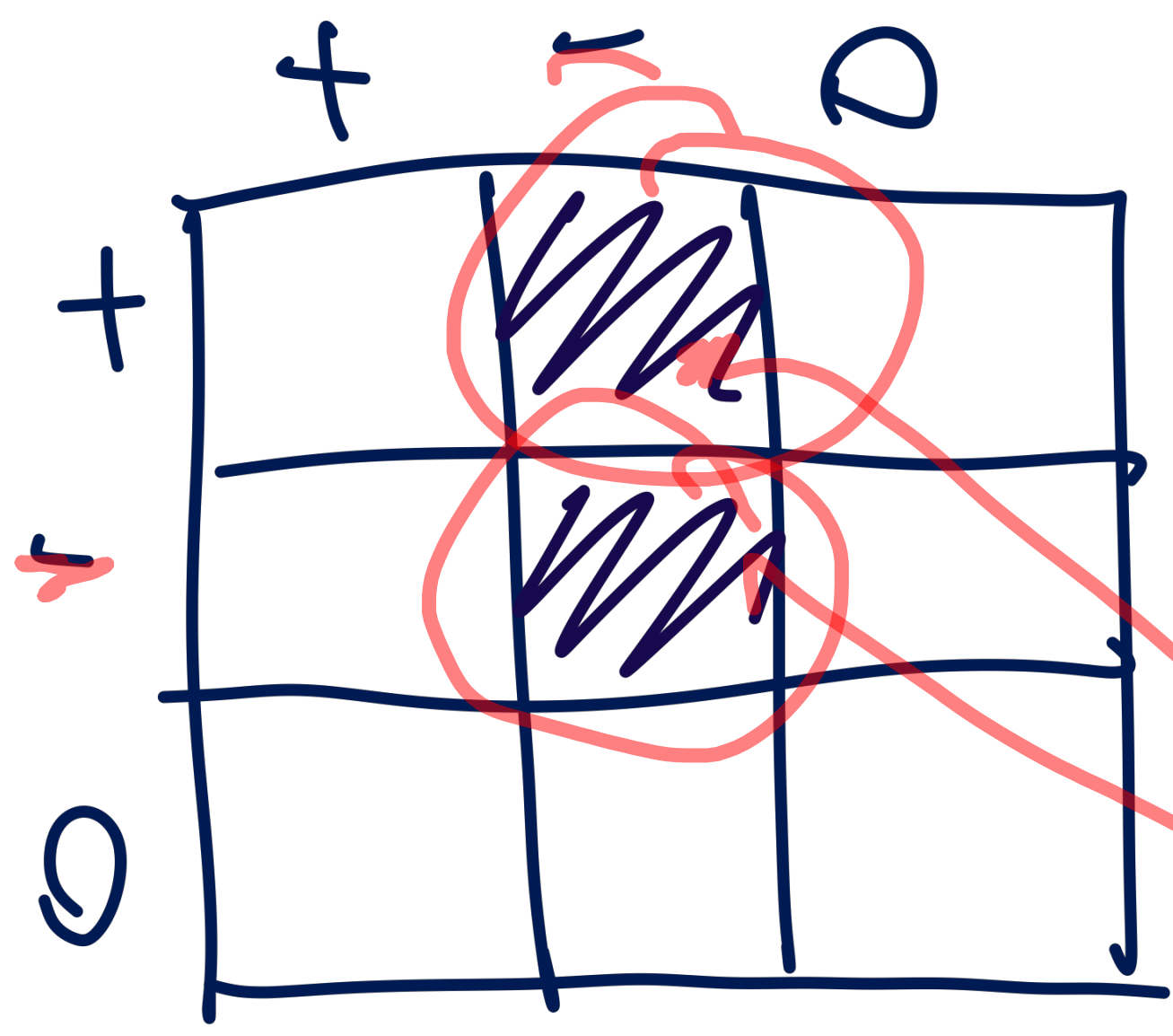
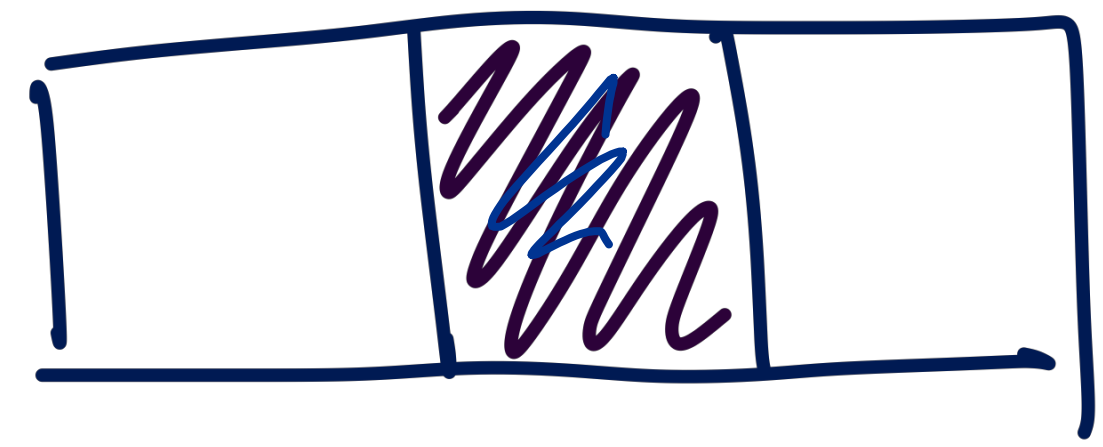
+ - 0



t=2

sentence 3

+ - 0



score  $(\begin{matrix} + & - & - \\ + & - & - \\ + & - & - \end{matrix})$

$$= \text{score}_0(+)+$$

$$\text{score}_1(-)+$$

$$\text{score}_2(-)+$$

$$+ \text{tscore}(+,-) + \text{tscore}(-,-)$$

$$q = 3 + 3 + 3 = 27$$

$$3 \times 3 \times 3$$

combinations = too slow

label sequence	unnormalized	probability
+ + +	-	softmax(-)
+ + -	-	
+ + 0	-	
+ - +	-	
+ - -	-	
+ - 0	-	
⋮	1	} sum = 1

# Conditional Random Fields

- Training ต้องใช้ algorithm ที่หา probability ได้เร็วๆ
- Decoding ต้องใช้ algorithm ที่หา label sequence ที่ดีที่สุดได้เร็วๆ  
หา sequence of labels ที่ดีที่สุด

# Decoding



ประโยค	เรื่องนี้คนแสดงนำหล่อ	แต่เนื้อเรื่องน่าเบื่อ	จบได้ جيدมาก
Sentiment	?	?	?

คำ	They	can	fish
POS tag	? Pron	? V	? V

# Training CRF

# Simple Classification

label	sentence
f	$s_1$   $A_1$
-	$s_2$
-	$s_3$
0	$s_4$   $A_2$
1	$s_5$
+	$s_6$   $A_3$
+	$s_7$
+	$s_8$

# Sequence Labeling

label sequence	sentence sequence
$[+, -, -]$	$[s_1, s_2, s_3]$
$[0, -]$	$[s_4, s_5]$
$[+, +, +]$	$[s_6, s_7, s_8]$

$3 \times 3 \times 3$  combinations = too slow

label sequence	unnormalized	probability
+ + +		
+ + 1		
+ + 0		
+ 1 +		
+ 1 1		
+ 1 0		
⋮		

sum = 1

# Objective Function

sequence of feature vectors

$$L(\theta) = - \sum \log P(\bar{Y} | \bar{X})$$

sequences  $\in D$

gradient

sequence of labels

[++0, --]

# Training Algorithm

- Forward-Backward algorithm
- Averaged Structured Perceptron

# Forward-Backward Algorithm

- คำนวณ log-likelihood ของ data (forward)
- คำนวณ gradient ของแต่ละ parameter (forward-backward)

# Averaged Structured Perceptron

minibatch

30 sequences

update

- Viterbi algorithm ในการ decode แล้วปรับแก้ parameter โดยไม่  
ต้องใช้ gradient



