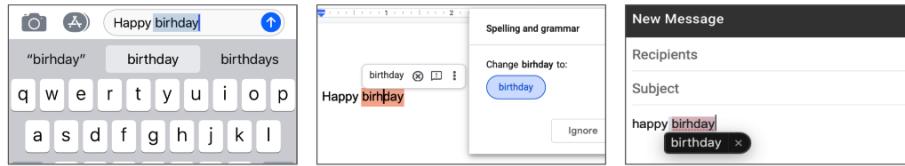


Autocorrect and Minimum Edit distance

• Autocorrect



To implement autocorrect in this week's assignment, you have to follow these steps:

- Identify a misspelled word
- Find strings n edit distance away: (these could be random strings)
- Filter candidates: (keep only the real words from the previous steps)
- Calculate word probabilities: (choose the word that is most likely to occur in that context)

• Building the model

① Identify the misspelled word

- 사건이 있다면 typo 일 수 있음

② Find Strings n edit distance away

- Edit: an operation performed on a string to change it
- Insert (add a letter) 'to': 'top', 'two' ...
- Delete (remove a letter) 'hat': 'ha', 'at', 'ht'
- Switch (swap 2 adjacent letters) 'eta': 'eat', 'tea'
- Replace (change 1 letter to another) 'jaw': 'jar', 'paw', ...

③ Filter candidates

- ②로 생성된 단어들 실제 단어만 취함

④ Calculate word probabilities

- language model을 이용해 수도 있지만 이전에는 word frequency 쓰기

Calculate word probabilities

Example: "I am happy because I am learning"

$$P(w) = \frac{C(w)}{V} \quad P(am) = \frac{C(am)}{V} = \frac{2}{7}$$

$P(w)$ Probability of a word

$C(w)$ Number of times the word appears

V Total size of the corpus

Word	Count
I	2
am	2
happy	1
because	1
learning	1

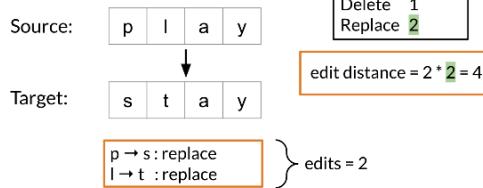
Total: 7

• Minimum edit distance

Minimum edit distance allows you to:

- Evaluate similarity between two strings
- Find the minimum number of edits between two strings
- Implement spelling correction, document similarity, machine translation, DNA sequencing, and more

Example:



Source: play → Target: stay

Cost: insert: 1, delete: 1, replace: 2

p → s

insert + delete: p → ps → s: 2
 delete + insert: p → # → s: 2
 replace: p → s: 2

	0	1	2	3	4
0	#	0	1		
1	p	1	2		
2					
3					
4					

Source: play → Target: stay

Cost: insert: 1, delete: 1, replace: 2

p → s

D[i, j] =

$$\min \begin{cases} D[i-1, j] + \text{del_cost} \\ D[i, j-1] + \text{ins_cost} \\ D[i-1, j-1] + \begin{cases} \text{rep_cost}; & \text{if } \text{src}[i] \neq \text{tar}[j] \\ 0; & \text{if } \text{src}[i] = \text{tar}[j] \end{cases} \end{cases}$$

	0	1	2	3	4
0	#	0	1		
1	p	1	2		
2	I	2			
3	a	3			
4	y	4			

Source: play → Target: stay

Cost: insert: 1, delete: 1, replace: 2

• Levenshtein distance

• Backtrace

• Dynamic programming

	0	1	2	3	4
0	#	0	1	2	3
1	p	1	2	3	4
2	I	2	3	4	5
3	a	3	4	5	4
4	y	4	5	6	5

Part of Speech Tagging

- POS Tagging: process of assigning a part of speech to a word.

Subject

- Markov chain
- Hidden Markov Models
- Viterbi algorithm

Part of speech tags:

lexical term	tag	example
noun	NN	something, nothing
verb	VB	learn, study
determiner	DT	the, a
wh-adverb	WRB	why, where
...	...	

Why not learn something?

WRB RB VB NN

Usage

- Identifying named entities
- Speech recognition
- Conference Resolution

POS 흐름 co-occurrence 를 추로시 사용가능 (LM과 연관화)

Markov Chains

- 다음 단어의 확률을 나타낼때 사용

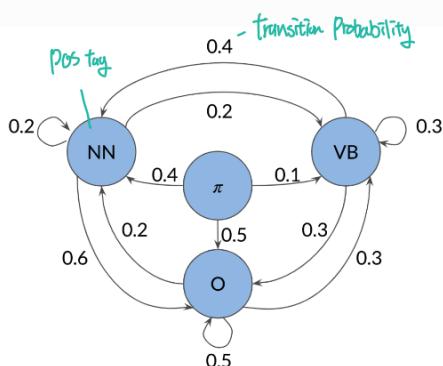
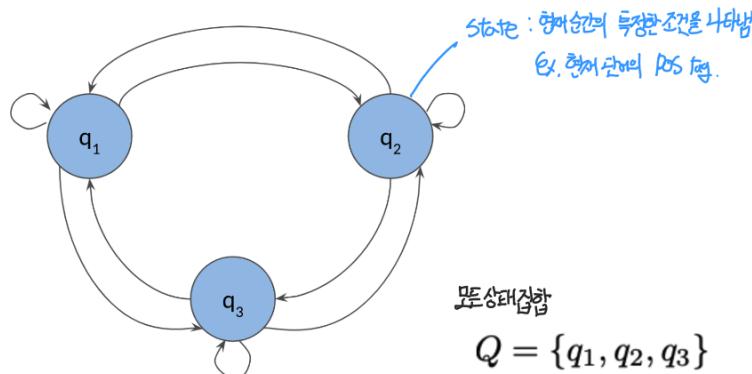
Why not learn swimming?

verb noun 확률↑

Why not learn swim?

verb verb 확률↓

- 확률 모델을 만드려면 POS tag와 word 확률이 필요함



Transition matrix

	NN	VB	O
π (initial)	0.4	0.1	0.5
NN (noun)	0.2	0.2	0.6
VB (verb)	0.4	0.3	0.3
O (other)	0.2	0.3	0.5

Initial distribution. Ex. 40% 확률로 문장이
NN으로 시작

Formal Notation

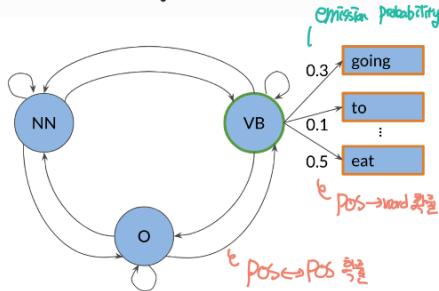
States

$$Q = \{q_1, \dots, q_N\}$$

Transition matrix

$$A = \begin{pmatrix} a_{1,1} & \dots & a_{1,N} \\ \vdots & \ddots & \vdots \\ a_{N+1,1} & \dots & a_{N+1,N} \end{pmatrix}$$

Hidden Markov Models • emission probability를 통해 State(POS)의 specific word의 확률을 나누는법



		going	to	eat	...
NN (noun)	0.5	0.1	0.02		
VB (verb)	0.3	0.1	0.5		
O (other)	0.3	0.5	0.68		

labelled dataset에 POS 태그에 대해 확률 계산

States

$$Q = \{q_1, \dots, q_N\}$$

Transition matrix

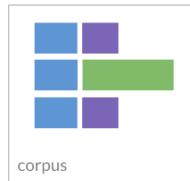
$$A = \begin{pmatrix} a_{1,1} & \dots & a_{1,N} \\ \vdots & \ddots & \vdots \\ a_{N+1,1} & \dots & a_{N+1,N} \end{pmatrix}$$

Emission matrix

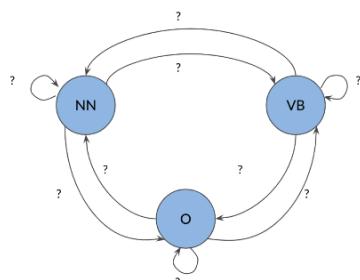
$$B = \begin{pmatrix} b_{11} & \dots & b_{1V} \\ \vdots & \ddots & \vdots \\ b_{N1} & \dots & b_{NV} \end{pmatrix}$$

$$\sum_{j=1}^V b_{ij} = 1$$

Calculating
Probabilities



transition probability: $\frac{\text{blue}}{\text{blue} + \text{purple}} = \frac{2}{3}$



1. Count occurrences of tag pairs $C(t_{i-1}, t_i)$: t_{i-1} $\stackrel{\text{def}}{=} \text{tag before}$ count

2. Calculate probabilities using the counts

$$P(t_i | t_{i-1}) = \frac{C(t_{i-1}, t_i)}{\sum_{j=1}^N C(t_{i-1}, t_j)}$$

t_{i-1} $\stackrel{\text{def}}{=} \text{anything}$ count

Populating
Transition
Matrix

	NN	VB	O
π	1	0	2
NN (noun)	0	0	6
VB (verb)	0	0	0
O (other)	6	0	8

① count $((t_1, t_2))$

<> in a station of the metro

<> the apparition of these faces in the crowd :

<> petals on a wet, black bough.

Ezra Pound - 1913

	NN	VB	O
π	1	0	2
NN	0	0	6
VB	0	0	0
O	6	0	8

② 확률 계산

$$P(\text{NN} | \text{O}) = \frac{C(\text{O}, \text{NN})}{\sum_{j=1}^N C(\text{O}, t_j)} = \frac{6}{14}$$

To generalize:

$$P(t_i | t_{i-1}) = \frac{C(t_{i-1}, t_i)}{\sum_{j=1}^N C(t_{i-1}, t_j)}$$

Smoothing

	NN	VB	O	
π	1+ ϵ	0+ ϵ	2+ ϵ	3+3* ϵ
NN	0+ ϵ	0+ ϵ	6+ ϵ	6+3* ϵ
VB	0+ ϵ	0+ ϵ	0+ ϵ	0+3* ϵ
O	6+ ϵ	0+ ϵ	8+ ϵ	14+3* ϵ

$$P(t_i | t_{i-1}) = \frac{C(t_{i-1}, t_i) + \epsilon}{\sum_{j=1}^N C(t_{i-1}, t_j) + N * \epsilon}$$

• 확률이 0이 되지 않도록

Populating Emission Matrix

Matrix

	in	a	...
NN (noun)	0		
VB (verb)	0		
O (other)	2		

<ss> in a station of the metro
 <ss> the apparition of these faces in the crowd:
 <ss> petals on a wet, black bough.

Ezra Pound - 1913

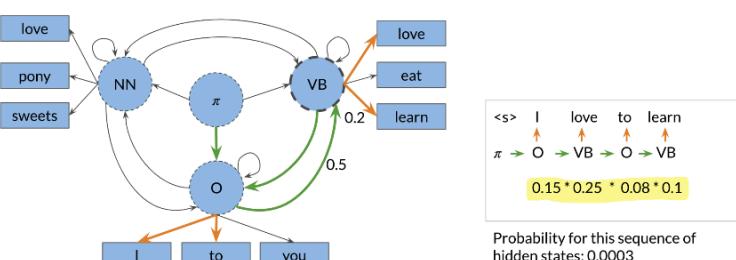
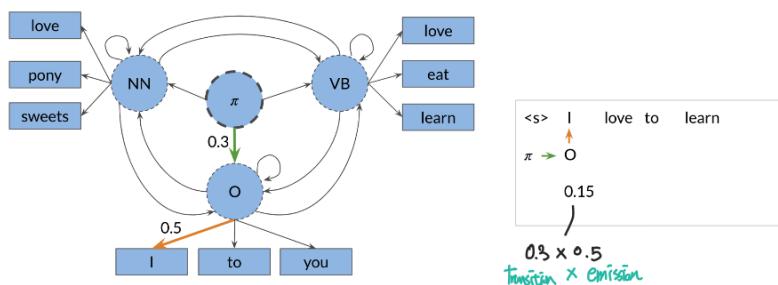
$$\begin{aligned}
 P(w_i | t_i) &= \frac{C(t_i, w_i) + \epsilon}{\sum_{j=1}^V C(t_i, w_j) + N * \epsilon} \\
 &= \frac{C(t_i, w_i) + \epsilon}{C(t_i) + N * \epsilon}
 \end{aligned}$$

Smoothing

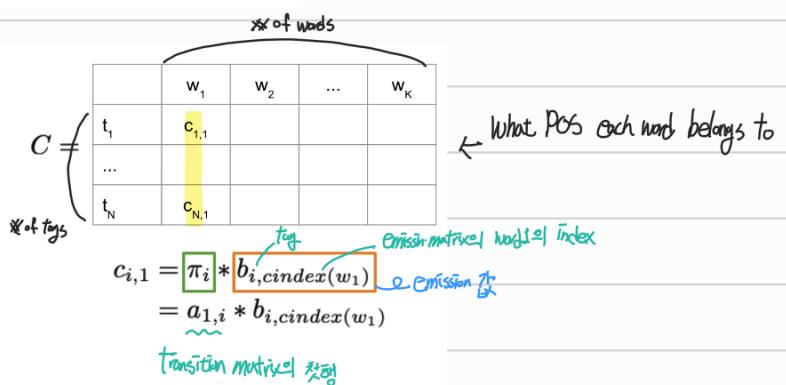
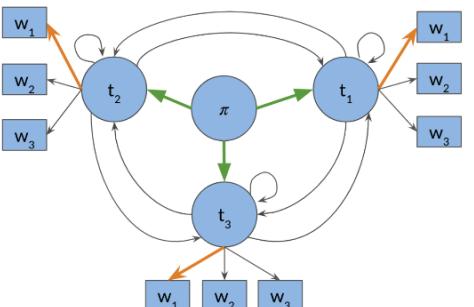
• $C(t_i, w_i)$: tag t_i 와 w_i 가 겹친 횟수 count

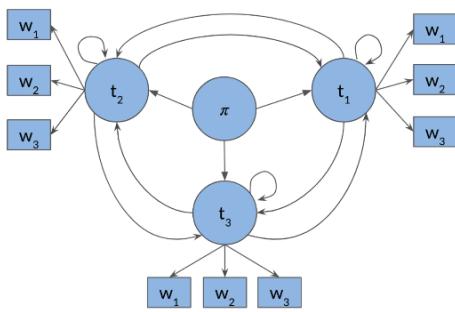
Viterbi : transition matrix + emission matrix로 확률 계산

The Viterbi Algorithm



Viterbi Initialization





	w ₁	w ₂	...	w _K
t ₁	d _{1,1}			
...				
t _N	d _{N,1}			

- What POS coming from

- w_1, \dots, w_N 의 sequence가 주제를 갖다

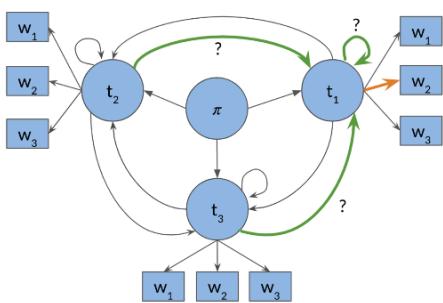
가장 높은 확률의 POS sequence를 찾을 때 사용

$$d_{i,1} = 0$$

처음에 다른 POS부터 시작해 0으로 초기화

Viterbi:

Forward Pass

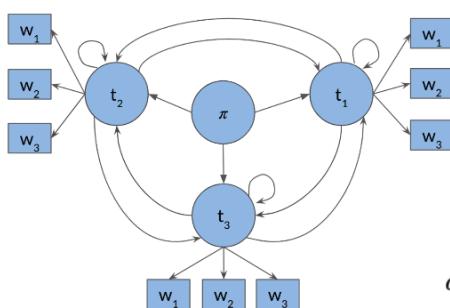


	w ₁	w ₂	...	w _K
t ₁	c _{1,1}	c _{1,2}		c _{1,K}
...				
t _N	c _{N,1}	c _{N,2}		c _{N,K}

$$c_{1,2} = \max_k c_{k,1} * a_{k,1} * b_{1,cindex(w_2)}$$

모든 [k] 대비 시 [2]의 최대값.

keep track of which was used to get max probability



	w ₁	w ₂	...	w _K
t ₁	d _{1,1}	d _{1,2}		d _{1,K}
...				
t _N	d _{N,1}	d _{N,2}		d _{N,K}

Probability

$$c_{i,j} = \max_k c_{k,j-1} * a_{k,i} * b_{i,cindex(w_j)}$$

$$d_{i,j} = \operatorname{argmax}_k c_{k,j-1} * a_{k,i} * b_{i,cindex(w_j)}$$

Index

Viterbi:

Backward Pass

$$C =$$

	w ₁	w ₂	w ₃	w ₄	w ₅
t ₁	0.25	0.125	0.025	0.0125	0.01
t ₂	0.1	0.025	0.05	0.01	0.003
t ₃	0.3	0.05	0.025	0.02	0.0000
t ₄	0.2	0.1	0.000	0.0025	0.0003

$$s = \operatorname{argmax}_i c_{i,K} = 1$$

마지막 col의 row index 시작!

$$D =$$

	w ₁	w ₂	w ₃	w ₄	w ₅
t ₁	0	1	3	2	3
t ₂	0	2	4	1	3
t ₃	0	2	4	1	4
t ₄	0	4	4	3	1

시작 col index

<s>	w1	w2	w3	w4	w5
π	$\leftarrow t_2$	$\leftarrow t_3$	$\leftarrow t_1$	$\leftarrow t_3$	$\leftarrow t_1$

∴ Reconstruct POS tags of sentence

Week 3

Autocomplete and Language Models

N-Grams Overview

- Calculate probabilities of certain words happening in a specific sequence.

<Application>

Speech recognition



$$P(\text{I saw a van}) > P(\text{eyes awe of an})$$

Spelling correction



"He entered the ship to buy some groceries" - "ship" a dictionary word

$$\bullet P(\text{entered the shop to buy}) > P(\text{entered the ship to buy})$$

Augmentative communication



Predict most likely word from menu for people unable to physically talk or sign.
(Newell et al., 1998)

N-grams and Probabilities

An N-gram is a sequence of N words

Corpus: I am happy because I am learning

Unigrams: { I, am, happy, because, learning }

Bigrams: { I am, am happy, happy because ... } ✖ I happy

Trigrams: { I am happy, am happy because, ... }

Corpus: This is great ... teacher drinks tea. $w_1 w_2 w_3 \dots w_{498} w_{499} w_{500}$ $m = 500$

• 가로상의

- $w_1^m = w_1 w_2 w_3 \dots w_m$
- $w_1^3 = w_1 w_2 w_3$
- $w_{m-2}^m = w_{m-2} w_{m-1} w_m$

Unigram

- Size of corpus $m = 7$.
- $P(I) = \frac{2}{7}$
- $P(happy) = \frac{1}{7}$

To generalize, the probability of a unigram is $P(w) = \frac{C(w)}{m}$

Bigram Probability:

Corpus: I am happy because I am learning

$$P(am|I) = \frac{C(I\ am)}{C(I)} = \frac{2}{2} = 1$$

$$P(happy|I) = \frac{C(I\ happy)}{C(I)} = \frac{0}{2} = 0 \quad \text{✖ I happy}$$

$$P(learning|am) = \frac{C(am\ learning)}{C(am)} = \frac{1}{2}$$

$$\text{Probability of a bigram: } P(y|x) = \frac{C(x\ y)}{\sum_w C(x\ w)} = \frac{C(x\ y)}{C(x)}$$

Trigram Probability:

To compute the probability of a trigram:

- $P(w_3 | w_1^2) = \frac{C(w_1^2 w_3)}{C(w_1^2)}$
- $C(w_1^2 w_3) = C(w_1 w_2 w_3) = C(w_1^3)$

N-gram Probability:

- $P(w_N | w_1^{N-1}) = \frac{C(w_1^{N-1} w_N)}{C(w_1^{N-1})}$
- $C(w_1^{N-1} w_N) = C(w_1^N)$

Sequence Probabilities • How to approximate N-gram probabilities?

(구하려면 sentence의 probability를 구해야함)

• 조건부 확률

- $P(B | A) = \frac{P(A, B)}{P(A)} \implies P(A, B) = P(A)P(B | A)$
- $P(A, B, C, D) = P(A)P(B | A)P(C | A, B)P(D | A, B, C)$

Ex.

$$P(\text{the teacher drinks tea}) = P(\text{the})P(\text{teacher the})P(\text{drinks the teacher})P(\text{tea} | \text{the teacher drinks})$$

• Corpus: 내가 구하려는 문장이 매우 드물게 등장한다. (특히 문장이 길 때) → 확률이 0에 가까움

• Markov Assumption: Only the last word matters

- Bigram $P(w_n | w_1^{n-1}) \approx P(w_n | w_{n-1})$ ← 시전 1개
- N-gram $P(w_n | w_1^{n-1}) \approx P(w_n | w_{n-N+1}^{n-1})$ ← 시전 N-1개

• 전체 문장 확률 구사

- $P(w_1^n) \approx \prod_{i=1}^n P(w_i | w_{i-1})$
- $P(w_1^n) \approx P(w_1) P(w_2 | w_1) \dots P(w_n | w_{n-1})$

Starting and
Ending Sentences

• N-gram language model 일 때 맨 앞에 N-1개 </s> 추가, 맨 끝에는 </s> 하나 추가

Corpus

< /s > Lyn drinks chocolate < /s >
< /s > John drinks tea < /s >
< /s > Lyn eats chocolate < /s >

$$P(\text{sentence}) = \frac{2}{3} * \frac{1}{2} * \frac{1}{2} * \frac{2}{2} = \frac{1}{6}$$

$$P(\text{John} | \text{</s>}) = \frac{1}{3}$$

$$P(\text{chocolate} | \text{eats}) = \frac{1}{2}$$

$$P(\text{</s>} | \text{tea}) = \frac{1}{1}$$

$$P(\text{Lyn} | \text{</s>}) = ? = \frac{2}{3}$$

The N-gram

- Count matrix

Language Model

- Bigram count matrix

"study I" bigram

Corpus: <s> I study I learn </s>

	<s>	</s>	I	study	learn
<s>	0	0	1	0	0
</s>	0	0	0	0	0
I	0	0	0	1	1
study	0	0	1	0	0
learn	0	1	0	0	0

Unique corpus N-grams.

자료는 bigram으로 되어야함.

Unique Corpus words

- Probability matrix로 변환

$$P(w_n | w_{n-N+1}^{n-1}) = \frac{C(w_{n-N+1}^{n-1}, w_n)}{C(w_{n-N+1}^{n-1})}$$

w_{n-1} row, col
Sum(row)

$$\text{sum}(row) = \sum_{w \in V} C(w_{n-N+1}^{n-1}, w) = C(w_{n-N+1}^{n-1})$$

- Probability matrix로 LM 생성 가능. 문장 확률 / 단어 예측 가능

To compute the probability of a sequence, you needed to compute:

$$P(w_1^n) \approx \prod_{i=1}^n P(w_i | w_{i-1})$$

To avoid underflow, you can multiply by the log.

$$\log(P(w_1^n)) \approx \sum_{i=1}^n \log(P(w_i | w_{i-1}))$$

Finally here is a summary to create the generative model:

Corpus:

<s> Lyn drinks chocolate </s>
<s> John drinks tea </s>
<s> Lyn eats chocolate </s>

- (<s>, Lyn) or (<s>, John)?
- (Lyn,eats) or (Lyn,drinks) ?
- (drinks,tea) or (drinks,chocolate)?
- (tea,</s>) - always

Algorithm:

- Choose sentence start
- Choose next bigram starting with previous word
- Continue until </s> is picked

Language Model

Evaluation

① Splitting the Data

Train/Val/Test splits

There are two main methods for splitting the data:

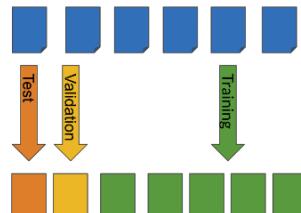
Smaller Corpora:

- 80% train
- 10% val
- 10% test

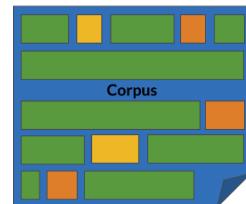
Larger Corpora:

- 98% train
- 1% val
- 1% test

- Continuous text



- Random short sequences



② Perplexity

- 사람이 쓴 문장: ↓, Random: ↑

$$PP(W) = P(s_1, s_2, \dots, s_m)^{-\frac{1}{m}}$$

$$PP(W) = \sqrt[m]{\prod_{i=1}^m \prod_{j=1}^{|s_i|} \frac{1}{P(w_j^{(i)} | w_{j-1}^{(i)})}}$$

$N_j^{(i)}$: i번째 문장의 j번째 단어

→ 모든 문장을 이으면 N 을 표현 가능

$$PP(W) = \sqrt[m]{\prod_{i=1}^m \frac{1}{P(w_i | w_{i-1})}}$$

To

$$\log PP(W) = -\frac{1}{m} \sum_{i=1}^m \log_2 (P(w_i | w_{i-1}))$$

Ex. speech recognition, question answering

OOV Words

- Vocabulary: set of unique words supported by language model
 - Closed: only fixed set of words
 - Open: encounter words outside the vocabulary

- Unknown word [문제]

- Create vocabulary V
- Replace any word in corpus and not in V by <UNK>
- Count the probabilities with <UNK> as with any other word

Corpus

<s> Lyn drinks chocolate </s>
<s> John drinks tea </s>
<s> Lyn eats chocolate </s>



Corpus

<s> Lyn drinks chocolate </s>
<s> <UNK> drinks <UNK> </s>
<s> Lyn <UNK> chocolate </s>

Min frequency f=2 → 이보다 높은 빈도수의 단어는 등록하지 않음

Vocabulary
Lyn, drinks, chocolate

Input query
<s> Adam drinks chocolate </s>
<s> <UNK> drinks chocolate </s>

Criteria to create the vocabulary

- Min word frequency f
- Max |V|, include words by frequency
- Use <UNK> sparingly (Why?)
- Perplexity - only compare LMs with the same V

Smoothing

$$P(w_n | w_{n-N+1}^{n-1}) = \frac{C(w_{n-N+1}^{n-1}, w_n)}{C(w_{n-N+1}^{n-1})} \text{ can be 0}$$

Hence we can add-1 smoothing as follows to fix that problem:

$$P(w_n | w_{n-1}) = \frac{C(w_{n-1}, w_n) + 1}{\sum_{w \in V} (C(w_{n-1}, w) + 1)} = \frac{C(w_{n-1}, w_n) + 1}{C(w_{n-1}) + V}$$

Add-k smoothing is very similar:

$$P(w_n | w_{n-1}) = \frac{C(w_{n-1}, w_n) + k}{\sum_{w \in V} (C(w_{n-1}, w) + k)} = \frac{C(w_{n-1}, w_n) + k}{C(w_{n-1}) + k * V}$$

back off

- If N-gram missing => use (N-1)-gram, ...: Using the lower level N-grams (i.e. (N-1)-gram, (N-2)-gram, down to unigram) distorts the probability distribution. Especially for smaller corpora, some probability needs to be discounted from higher level N-grams to use it for lower level N-grams.
- Probability discounting e.g. Katz backoff: makes use of discounting.
- "Stupid" backoff: If the higher order N-gram probability is missing, the lower order N-gram probability is used, just multiplied by a constant. A constant of about 0.4 was experimentally shown to work well.

back off example

Corpus

<s> Lyn drinks chocolate </s>
 <s> John drinks tea </s>
 <s> Lyn eats chocolate </s>

$$P(\text{chocolate} | \text{John drinks}) = ?$$

$$0.4 \times P(\text{chocolate} | \text{drinks})$$



• Interpolation

$$\hat{P}(w_n | w_{n-2} w_{n-1}) = \lambda_1 \times P(w_n | w_{n-2} w_{n-1}) + \lambda_2 \times P(w_n | w_{n-1}) + \lambda_3 \times P(w_n)$$

Where

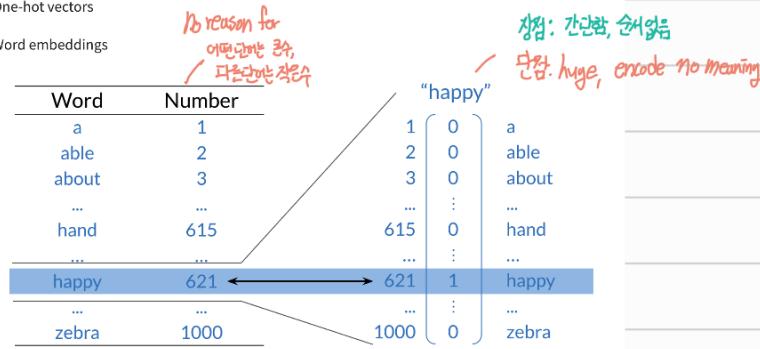
$$\sum_i \lambda_i = 1$$

Week 4

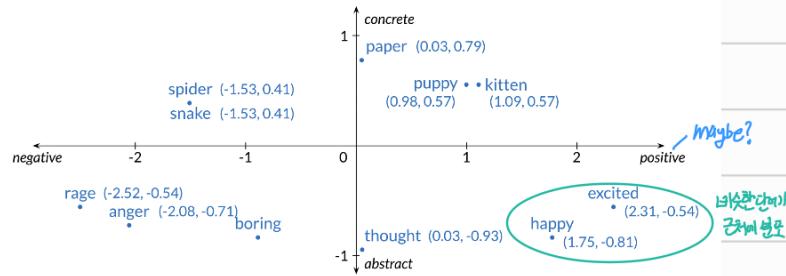
Word Embeddings

• Basic Word Representations

- Integers
- One-hot vectors
- Word embeddings



• Word Embeddings



장점: Low dimensions (less than V)

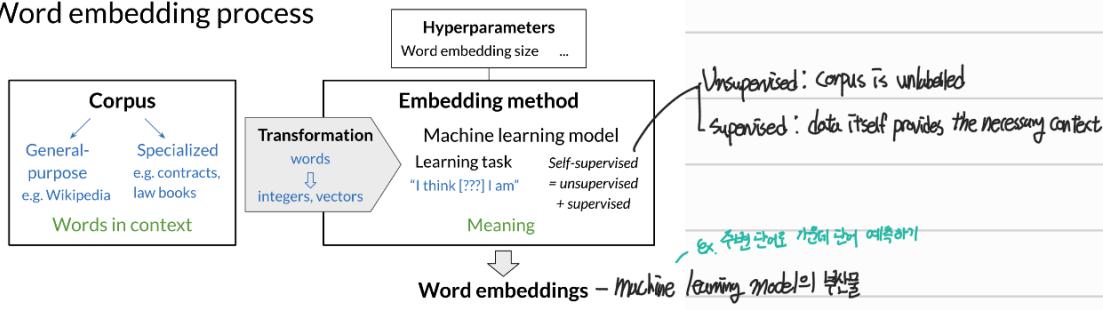
Encode Meaning

• Coordinate Word聯繫에 어떤 성분을 주고 생각

How to create
Word Embeddings

- Need: ① Corpus ② embedding method
- Context: 각 word embedding이 의미를 부여함

Word embedding process



Word Embedding

Methods

Classical Methods

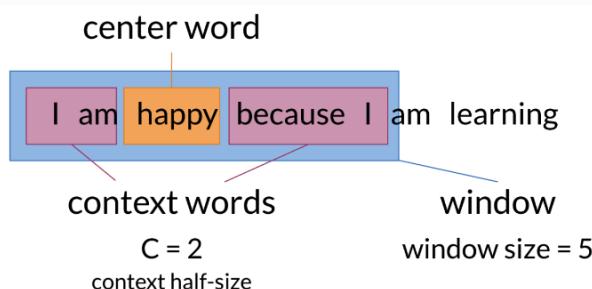
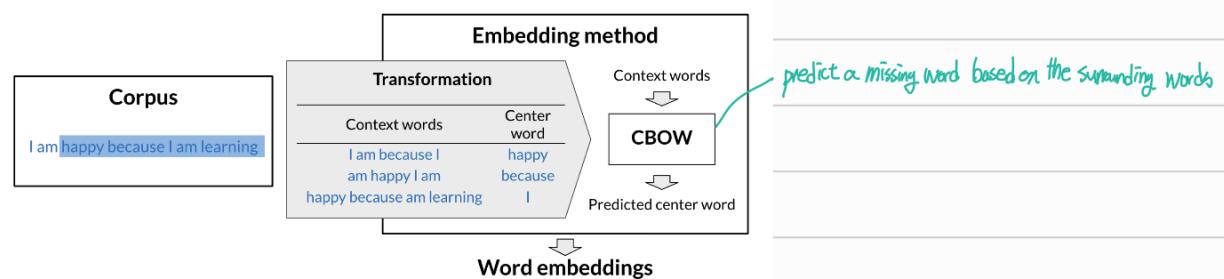
- word2vec (Google, 2013)
- Continuous bag-of-words (CBOW)*: the model learns to predict the center word given some context words.
- Continuous skip-gram / Skip-gram with negative sampling (SGNS)*: the model learns to predict the words surrounding a given input word.
- Global Vectors (GloVe) (Stanford, 2014)*: factorizes the logarithm of the corpus's word co-occurrence matrix, similar to the count matrix you've used before.
- fastText (Facebook, 2016)*: based on the skip-gram model and takes into account the structure of words by representing words as an n-gram of characters. It supports out-of-vocabulary (OOV) words.

Deep learning, contextual embeddings

In these more advanced models, words have different embeddings depending on their context. You can download pre-trained embeddings for the following models.

- BERT (Google, 2018)
- ELMo (Allen Institute for AI, 2018)
- GPT-2 (OpenAI, 2018)

CBOW model



Cleaning and Tokenization

- NLP algorithm implement 전에 거친 절차

- Letter case "The" == "the" == "THE" → lowercase / uppercase
- Punctuation , ! . ? → . " ‘ ‘ ’ ’ ” → ∅ ... !! ??? → ..
- Numbers 1 2 3 5 8 → ∅ 3.14159 90210 → as is / <NUMBER>
- Special characters ⚡ \$ € § ¶ ** → ∅
- Special words 😊 #nlp → :happy: #nlp

Transforming Words

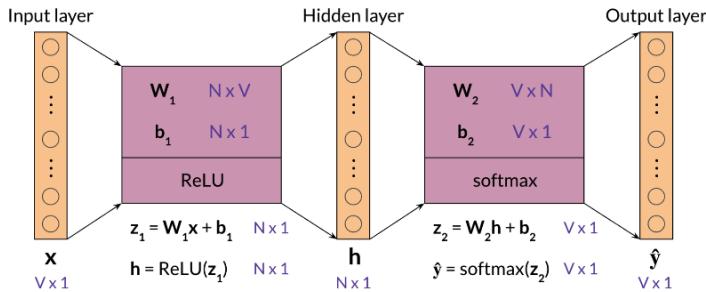
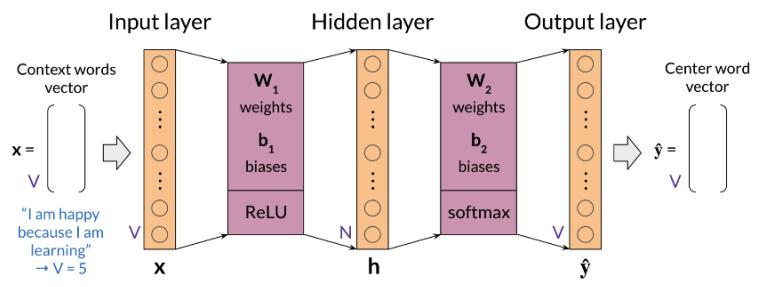
into Vectors

$$\left\{ \begin{array}{c} \text{I} \\ \text{am} \\ \text{because} \\ \text{happy} \\ \text{I} \\ \text{learning} \end{array} \right\} + \left\{ \begin{array}{c} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{array} \right\} + \left\{ \begin{array}{c} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{array} \right\} + \left\{ \begin{array}{c} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{array} \right\} \right\} / 4 = \left\{ \begin{array}{c} 0.25 \\ 0.25 \\ 0 \\ 0.5 \\ 0 \end{array} \right\}$$

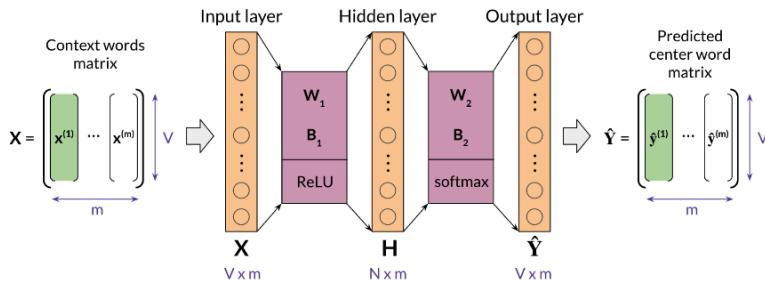
- One-hot vector 사용

Context words	Context words vector	Center word	Center word vector
I am because I	[0.25; 0.25; 0; 0.5; 0]	happy	[0; 0; 1; 0; 0]

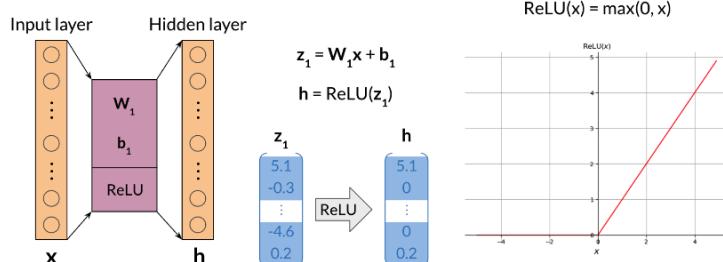
Architecture for the CBOW Model



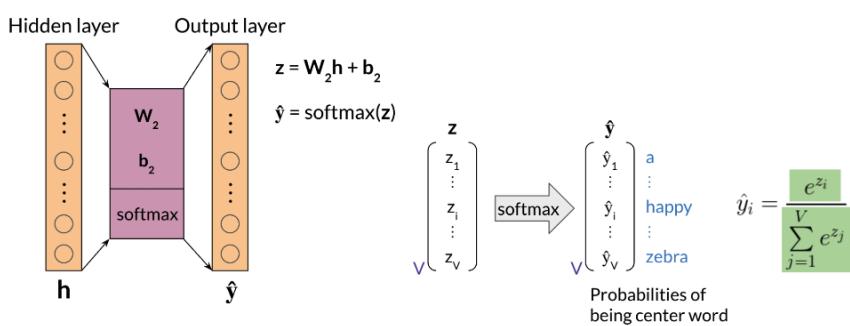
Batch Input dimension



Activation Functions - ReLU



Softmax : Vector → probability distribution



Training a CBOW model

- Cost function

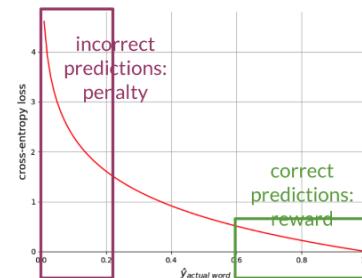
• Cross-entropy Error J_{CE}

$$J = - \sum_{k=1}^V y_k \log \hat{y}_k$$

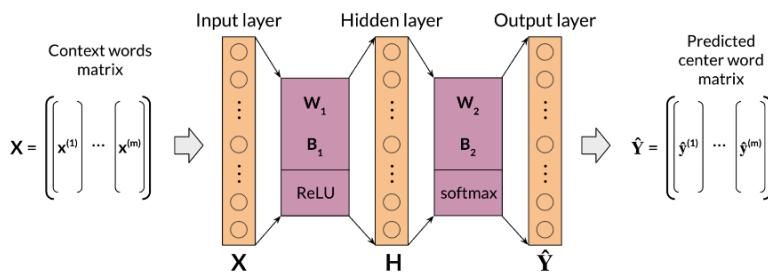
$$J = -\log \hat{y}_{\text{actual word}}$$

y		\hat{y}	
0	am	0.96	
0	because	0.01	
1	happy	0.01	
0	I	0.01	
0	learning	0.01	

$$\rightarrow J = 4.61$$



Forward Propagation



$$Z_1 = W_1 X + B_1$$

$$H = \text{ReLU}(Z_1)$$

$$Z_2 = W_2 H + B_2$$

$$\hat{Y} = \text{softmax}(Z_2)$$

$$J_{batch} = -\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^V y_j^{(i)} \log \hat{y}_j^{(i)}$$

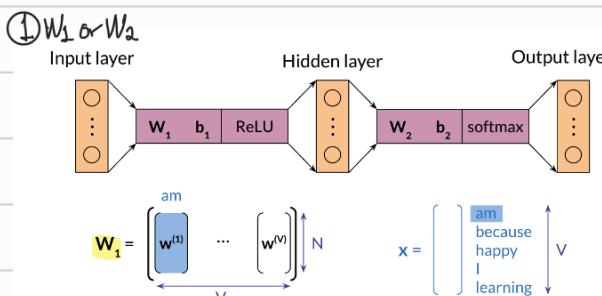
Backpropagation and Gradient descent

$$\begin{aligned} \mathbf{W}_1 &:= \mathbf{W}_1 - \alpha \frac{\partial J_{batch}}{\partial \mathbf{W}_1} \\ \mathbf{W}_2 &:= \mathbf{W}_2 - \alpha \frac{\partial J_{batch}}{\partial \mathbf{W}_2} \\ \mathbf{b}_1 &:= \mathbf{b}_1 - \alpha \frac{\partial J_{batch}}{\partial \mathbf{b}_1} \\ \mathbf{b}_2 &:= \mathbf{b}_2 - \alpha \frac{\partial J_{batch}}{\partial \mathbf{b}_2} \end{aligned}$$

learning rate

Calculate partial derivatives w.r.t weights & biases.

Extracting Word Embedding Vectors



If you were to use w_1 , each column will correspond to the embeddings of a specific word. You can also use w_2 as follows:

$$\mathbf{W}_2 = \begin{pmatrix} \mathbf{w}_2^{(1)} \\ \vdots \\ \mathbf{w}_2^{(V)} \end{pmatrix}$$

$$\mathbf{x} = \begin{pmatrix} \mathbf{am} \\ \mathbf{because} \\ \mathbf{happy} \\ \mathbf{I} \\ \mathbf{learning} \end{pmatrix}$$

② \mathbf{W}_1 과 \mathbf{W}_2 의 평균 취하기

$$\mathbf{W}_3 = 0.5 (\mathbf{W}_1 + \mathbf{W}_2^T) = \begin{pmatrix} \mathbf{w}_3^{(1)} \\ \vdots \\ \mathbf{w}_3^{(V)} \end{pmatrix}$$

Evaluating Word Embeddings:

- Test relationships between words

Ex. Semantic analogies: France : Paris = Italy : ?

Syntactic analogies: Seen : saw = been : ?

Intrinsic Evaluation

Test relationships between words

- Analogies
- Clustering
- Visualization



Extrinsic Evaluation

- tests embeddings on external tasks like named entity recognition, parts-of-speech tagging, etc..

- + Evaluates actual usefulness of embeddings
- - Time Consuming
- - More difficult to trouble shoot