

Lect 01- Intro

NLP

understand language

Communication

Human - Machine: QA, dialog

Human - Human: MT, spellcheck



NLP Task

Method of creating NLP

1. Rules: Bag of words
2. Prompting: Prompt w/o train
3. Finetune: learned from $\langle X, Y \rangle$

Input X
Text
Text
Text
Text
Image

Output Y
Corpus Text
Text in Other Language
Text Label
Linguistic Structure

Task
Language Modeling
Translation
Text Classification
Language Analysis
Image Captioning

Best Integrity (Thai)
split word if meaning not change
 $\text{ແມ່ນ} \neq \text{ແມ່ນ} + \text{ນັ້ນ}$
 $\text{ກວມພັນ} \sim \text{ກວມ} + \text{ພັນ}$

Tokenization

1. Word - High OOV
2. Subword - moderate
3. Char - no meaning

Tokenizer Algo

1. Longest match - match current word long as possible in Vocab

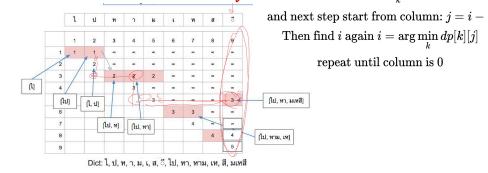
- ປຳເນົດນິກາ Start scanning with “ i ” as the starting point
- ປຳເນົດນິກາ Keep scanning ...
- ປຳເນົດນິກາ No more words start with “ ປຳເນົດ ”, move to the next point
- ປຳເນົດນິກາ

2. Maximal match (DP O(n^2))

$\text{dp}[i][j] = \text{fewest words from } [1, j] (\text{all in vocab})$, and current word start from $[i, j]$

$$dp_{i,j} = \begin{cases} \min_{k \leq i-1} dp_{k,i-1} + 1, & \text{if } sent_{i,j} \in V \\ \infty, & \text{otherwise} \end{cases}$$

basecase: $dp[i][1], i > 0$



ML => better to handle OOV

One hot: sparse = curse of dim, no meaning

neural embedding: learnable mapping table

CLASS torch.nlp.Embedding(num_embeddings, embedding_dim, padding_idx=None, max_norm=None, norm_type=2.0, scale_grad_by_freq=False, sparse=False, weight=None, freeze=False, device=None, dtype=None) [SOURCE]

embedding is actually a OHE then feed into linear layer

Lecture 02 - Language Model

LM - probabilistic model to find probability that an input sentence will occur

$$\text{Full estimation: } p(w_{1:n}) = p(w_1) \prod_{i=1}^n p(w_i | w_{1:i-1})$$

Rule-based LM

$$\text{N-gram: } Pr(w_i | w_{1:i-1}) = Pr(w_i | w_{i-n+1:i-1}) = \frac{c(w_{i-n+1:i})}{c(w_{i-n+1:i-1})} / Pr(< s >) = 1$$

Example: Bi-gram

• Estimating Bigram Probability

Uni-gram counting	
i	2533
am	927
to	2417
eat	746
chinese	158
food	1093
lunch	341
spend	278

Bi-grams counting (columns given row)	
i	= $c(w_{i-1:i}) = c(w_{1:2}) = 2533$
am	= $c(w_{2:3}) = c(w_{1:2}) = 927$
to	= $c(w_{3:4}) = c(w_{2:3}) = 2417$
eat	= $c(w_{4:5}) = c(w_{3:4}) = 746$
chinese	= $c(w_{5:6}) = c(w_{4:5}) = 158$
food	= $c(w_{6:7}) = c(w_{5:6}) = 1093$
lunch	= $c(w_{7:8}) = c(w_{6:7}) = 341$
spend	= $c(w_{8:9}) = c(w_{7:8}) = 278$

Evaluation

Extrinsic: downstream task - require human

Intrinsic: **Perplexity** ---***

$$\text{Perplexity}(P) = \sqrt[N]{\prod_{i=1}^N P(w_i)} = 2^{H(P)}, H(P) = -\frac{1}{N} \sum_{i=1}^N \log_2 P(w_i)$$

N-gram problem

1. Unknown words: OOVs => word that is not found in train set

handle ① Assign to normal word \rightarrow edit approach

• lower freq. than Min threshold \rightarrow treat as unk (replace on training set)

• top n freq.

e.g. training data has 1200 distinct words

\rightarrow use top 1000 highest frequency as a vocab

\rightarrow the rest (200) replace to unk tokens

\rightarrow and set $p(\text{unk}) = \frac{n(\text{unk})}{1200} = \frac{200}{1200}$

② set $p(\text{unk}) = \frac{1}{1200} = \frac{1}{1200}$ (from above example)

2. Zero: $\Rightarrow n\text{-gram count} = 0 \Rightarrow$ Perplexity is divided by 0

1. Laplace Smoothing

↳ add 1 to all words count.

$$P(w_i) = \frac{c_i}{N} \xrightarrow{\text{add 1}} P_{\text{Laplace}}(w_i) = \frac{c_i + 1}{N + V} \quad c_i = \text{ngram count}$$

word	freq	count	prob
want	2	0	0
am	0	0	0
to	1	0	0
eat	0	0	0
chinese	1	0	0
food	15	0	0
lunch	2	0	0
spend	1	0	0

2. LM Interpolation

↳ example: for tri-gram there is no $w_{n-2}w_{n-1}w_n$

↳ use bigram to estimate $P(w_n | w_{n-1})$ instead

$$\hat{P}(w_n | w_{n-2}w_{n-1}) = \lambda_1 P(w_n) + \lambda_2 P(w_n | w_{n-1}) + \lambda_3 P(w_n | w_{n-2}w_{n-1}) + \lambda_4 C$$

3. Stupid Backoff

back to n-1 gram until count != 0

4. Kneser-Ney Smoothing

backoff to unigram not making sense

d = abs discount (\Rightarrow count in train & test not equal) \Rightarrow penalty = 0.75

- If our corpus contains these bigrams
- $\text{Sam Francisco} = 1/4 = 0.25$
- $\text{glasses} = (3/4) = 0.75$
- Now, a word like “Francisco” will have low P continuation

type $(., w)$ = unique bigram count that end $w \Rightarrow$ fix col sum row

type $(w, .)$ = unique bigram count that start $w \Rightarrow$ fix row sum col

type $(., .)$ = all unique bigram count

$P_{\text{cont}}(w) = \frac{\text{type}(., w)}{\text{type}(., .)}$ all bigram end with w over all unique bigram

$\lambda(w_{i-1}) = \frac{d}{\text{cunigrams}(w_{i-1})} \text{type}(w_{i-1}, .)$, weight to make $\sum_{w \in W} P_{\text{cont}}(w, w_{i-1}) = 1$

$P_{\text{train}}(w_i | w_{i-1}) = \max(c_{\text{bigram}}(w_{i-1}, w_i), d, 0) + \lambda(w_{i-1}) P_{\text{cont}}(w, w_{i-1}) / \text{cunigrams}(w_{i-1})$

+ Example: a bigram Kneser-ney (cont.)

Compute the log-likelihood of the sentence “ am Sam Sam am ”

$$Pr(\text{am} | \text{am}) = \max(0.75, 0) / (0.75 + 0.25) = 0.6666666666666667$$

$$Pr(\text{am} | \text{Sam}) = \max(2 * 0.75, 0) / (0.75 * 2 + 0.25) = 0.8078$$

$$Pr(\text{Sam} | \text{am}) = \max(2 * 0.75, 0) / (0.75 * 2 + 0.25) = 0.8078$$

$$Pr(\text{Sam} | \text{Sam}) = \max(2 * 0.75, 0) / (0.75 * 2 + 0.25) = 0.8078$$

$$LL = \ln(0.34049) + \ln(0.8078) + \ln(0.8078) = -4.6892$$

↳ type $(., w)$ = $\text{type}(., w) + 1$

↳ type $(w, .)$ = $\text{type}(w, .) + 1$

↳ type $(., .)$ = $\text{type}(., .) + 2$

↳ sum column = type $(., .)$

↳ sum exists, t or o

Neural LM

Cost function

$$\mathcal{J} = -\frac{1}{T} \sum_{t=1}^T y_t \log \hat{p}_t \quad \rightarrow \text{cost only matched label}$$

Fact: \hat{p}_t is Perplexity \Rightarrow our probabilistic model tried to minimize perplexity

RNN based - RNN, GRU, LSTM

Lecture 02 - Subword (BPE, Wordpiece, Unigram, Sentencepiece)

Byte-Pair encoding (BPE) \rightsquigarrow Huffman Code like \Rightarrow merge most frequent token

Corpus: (“hug”, 10), (“pig”, 5), (“pun”, 12), (“bun”, 4), (“hugs”, 5)

epoch 1: hu = 10 + 5, pu = 5 + 12, bu = 12 + 4, hg = 5, ug = 10 + 5 + 5

↳ most frequent merge in corpus

usage:

↳ trained 3 epochs merging rule

How to use (test data)

- bug = “b”, “ug” (“b” in dict)

- mug = “m”, “ug” (“m” not in dict)

- thug = “[UNK”, “ug” (“t” not in dict)]

Wordpiece same alg. as BPE but we score instead of frequency

• NP score = $\frac{1}{2}(\text{char}_1, \text{char}_2) + \frac{1}{2}(\text{char}_1, \text{char}_3)$ \rightsquigarrow word like correlation (p)

• NP score = $\frac{1}{2}(\text{char}_1, \text{char}_2) + \frac{1}{2}(\text{char}_2, \text{char}_3)$ \rightsquigarrow key diff. NP regards “start of word”

↳ bug = (h, #H, #H) \rightsquigarrow bug = “b”

↳ mug = (m, #M, #M) \rightsquigarrow mug = “m”

↳ thug = ([UNK], #H, #H) \rightsquigarrow thug = “[UNK”]

↳ How to use “the longest subword”

- bug = “[UNK”, “#H”]

- If not possible to find subwords, tokenize the whole word as UNK

- bug = “[UNK”]

- bug = “[UNK”] \rightsquigarrow bug = “[UNK”]

↳ How to use “the longest subword”

- bug = “[UNK”, “#H”]

- If not possible to find subwords, tokenize the whole word as UNK

- bug = “[UNK”]

↳ How to use “the longest subword”

- bug = “[UNK”, “#H”]

- If not possible to find subwords, tokenize the whole word as UNK

- bug = “[UNK”]

↳ How to use “the longest subword”

- bug = “[UNK”, “#H”]

- If not possible to find subwords, tokenize the whole word as UNK

- bug = “[UNK”]

↳ How to use “the longest subword”

- bug = “[UNK”, “#H”]

- If not possible to find subwords, tokenize the whole word as UNK

- bug = “[UNK”]

↳ How to use “the longest subword”

- bug = “[UNK”, “#H”]

- If not possible to find subwords, tokenize the whole word as UNK

- bug = “[UNK”]

↳ How to use “the longest subword”

- bug = “[UNK”, “#H”]

- If not possible to find subwords, tokenize the whole word as UNK

- bug = “[UNK”]

↳ How to use “the longest subword”

- bug = “[UNK”, “#H”]

- If not possible to find subwords, tokenize the whole word as UNK

- bug = “[UNK”]

↳ How to use “the longest subword”

- bug = “[UNK”, “#H”]

- If not possible to find subwords, tokenize the whole word as UNK

- bug = “[UNK”]

↳ How to use “the longest subword”

- bug = “[UNK”, “#H”]

- If not possible to find subwords, tokenize the whole word as UNK

- bug = “[UNK”]

↳ How to use “the longest subword”

- bug = “[UNK”, “#H”]

- If not possible to find subwords, tokenize the whole word as UNK

- bug = “[UNK”]

↳ How to use “the longest subword”

- bug = “[UNK”, “#H”]

- If not possible to find subwords, tokenize the whole word as UNK

- bug = “[UNK”]

↳ How to use “the longest subword”

- bug = “[UNK”, “#H”]

- If not possible to find subwords, tokenize the whole word as UNK

- bug = “[UNK”]

↳ How to use “the longest subword”

- bug = “[UNK”, “#H”]

- If not possible to find subwords, tokenize the whole word as UNK

- bug = “[UNK”]

↳ How to use “the longest subword”

- bug = “[UNK”, “#H”]

- If not possible to find subwords, tokenize the whole word as UNK

- bug = “[UNK”]

↳ How to use “the longest subword”

- bug = “[UNK”, “#H”]

- If not possible to find subwords, tokenize the whole word as UNK

- bug = “[UNK”]

↳ How to use “the longest subword”

- bug = “[UNK”, “#H”]

- If not possible to find subwords, tokenize the whole word as UNK

- bug = “[UNK”]

positional encoding => same word = same embedding

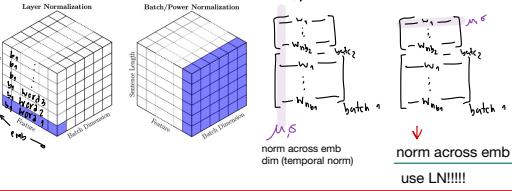
$$PE_{(pos,2i)} = \sin(pos/10000^{2i/dim})$$

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/dim})$$

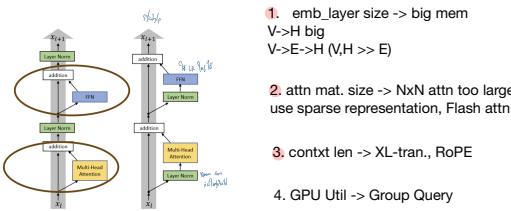
add positional information

type of PE: 1. absolute(fix context len) 2. relative

layer norm



More transformer



5. PreLN - solve gradient problem.

L04 - Token clf

e.g., PoS Tag

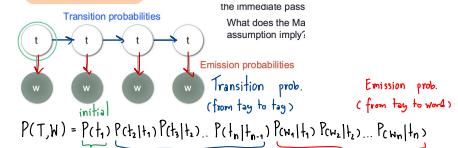
Traditional method

1. Seq method

$$T = \text{list of tag}, W = \text{list of word}, |T|=|W|=n$$

Objective $\arg \max P(T|W)$, i.e. find tags for each pos.

Modeling $P(T, W)$ instead (generative) w/ Markov Assumption



How to search for T that maximize $P(T, W)$ => Viterbi Algor (DP)

Define $\Pi[i, t]$ as the maximum log probability of the tag sequence from position 1 to i with $t_i = t$.

Initialization:

$$\Pi[0, < s >] = 0$$

$$\Pi[0, t] = -\infty \quad \forall t \neq < s >$$

Recursion:

$$\Pi[i, t] = \max_{t' \in T} \left(\Pi[i-1, t'] + \underbrace{\log p(t|t')}_\text{transition} + \underbrace{\log p(w_i|t')}_\text{emission} \right)$$

$$\text{backpointer}[i, t] = \arg \max_{t' \in T} (\Pi[i-1, t'] + \log p(t|t') + \log p(w_i|t'))$$

example:

	To N	To NN
From N	0.6	0.4
From NN	0.5	0.5

	I	eat	Chinese
State N	0.8	0.01	0.19
State NN	0.1	0.45	0.45

Recursive Step:

$$\pi[i, t] = \max_{t'} (\pi[i-1, t'] + \log p(t|t') + \log p(w_i|t'))$$

Not guarantee max

Viterbi Problem => Too slow O(T^2W) => use beam search is fine

Modeling $P(T|W)$ directly : Condition Random Field (CRF)

True discriminative => better for classification
add dependencies across all words



$$\text{Model: } P(T|W) = \prod_t P(t_i | t_{i-1}, \dots, t_1, W)$$

But W is too large => We calculate it via logit and softmax instead

$\text{softmax}(w_i | t_i)$ function ...

NN of PoS => BERT based, use NN w/ CRF

L05 - Word Representation

Symbolic - based on some rule(human labor) e.g. wordnet

Distributed: computed from word around it

sparse

- Term Frequency (Raw frequency) - doc info retrieval
- Co-occurrence -> small context (windowing)

3. Positive pointwise mutual info (PPMI)

$$\text{PMI}(w, c) = \log_2 \frac{P(w, c)}{P(w)P(c)}$$

target word context word

occur together occur independently

measure informative? -> PMI

$$\text{PPMI}(w, c) = \max(\text{PMI}(w, c), 0)$$

independent => not cont

4. TF-IDF

!designed for only docs similarity (not for word)

$$TF(w, \text{doc}) = \frac{\text{freq}(w \text{ in doc})}{\text{wordcount(doc)}}$$

$$IDF(w) = \log \frac{\# \text{docs}}{\#\text{docs that contain } w}$$

$$TF - IDF(w, \text{doc}) = TF(w, \text{doc}) - IDF(w)$$

IDF occur so frequent
= no sig. meaning

TF-idf variant: BM-25 => designed for RAG dense

handle typos + capture synonyms

How 2 train = train on target task, LM task, supervised task

Simple NN model

CBOV: give context predict target word -> just avg embds

Skip-gram: give target word predict context

New idea => negative sampling => idea to contrastive learn.

use some sample as loss (instead of whole vocab)

$o = \text{otherword}$ $c = \text{center(target) word}$

$$J_t(\theta) = \log \sigma(u_o^T v_c) + \sum_{i=1}^k \mathbb{E}_{j \sim P(w)} [\log \sigma(-u_j^T v_c)]$$

same cls diff class (negative)

Context word (positive, +1) Negative samples (negative, -1)

temperature(ing) (first era)

$$P^{\frac{1}{2}}(w) = \frac{(w)^{\frac{1}{2}}}{\sum w'(w')^{\frac{1}{2}}}$$

Spread out prob to rare word like languages

Pretrained-Word2Vec => fastText ~ handle type

add character n-gram in embedding

$$\text{apple} \Rightarrow f(\text{unk} \rightarrow) + w_1 f(\text{app} \rightarrow) \dots \approx \text{apple}$$

$$\text{orange} \Rightarrow f(\text{unk} \rightarrow) + w_1 f(\text{ow} \rightarrow) + \dots \approx \text{orange}$$

if not add, they are treat as same embed => not good

L05 - Text representation

Latent(emb) space -> lower dim representation

embd -> trained on supervised or self-supervised

Contrastive learning - similar thing = similar embedding (close angle)

Triplet loss -> pick anchor, positive, negative

$$* \sum_i^N [\|f(x_i^a) - f(x_i^p)\|^2 - \|f(x_i^a) - f(x_i^n)\|^2 + d]$$

positive negative margin

$f(x)$ is expensive compute all pair

InfoNCE -> 1 anchor, 1 positive, many negative

$$\text{InfoNCE Loss}_a = \frac{-1}{B \cdot P} \sum_{i=1}^B \sum_{p \in P(i)} \frac{e^{-\text{cost}(f(p))}}{\sum_{n \in N(i)} e^{-\text{cost}(f(n))}}$$

cost similarity

soft score

Might also use temperature here

negative mining => choose only hard case

MUSE - sentence encoding

DAN = average context + feed to MLP

USE - universal sentence encoder

Skip-thought(sentence level skip-gram), res-pred, NLI

use attention or DAN to emb sent. Then using pretrained (multitask)

MUSE - Multilingual USE

Bert-based embedding



Bert as sent. representation

- use NSP(sent. summary emb)
- use Pooled embedded emb on all task => shared encoder

Consert => contrastive self-supervised

the positive embedding = augmented embedding, other is negative

SimCSE=> from Consert apply only dropout is sufficient

L06 - Text Classification

task: sentiment, authorship id, spam filtering

Set of docs: D Set of classes: C

BOW(Naive Bayes) -> no positional (just wc)

count only buzzword (human defined)

$$\arg \max_C P(C|D) = \arg \max_C P(D|C)P(C)$$

(generative = MAP) $P(x) = \text{how many } x \text{ occurs in all docs}$

sampledD hair color eye color weight apply sum sum

S1 black dark overweight no no $x_1 = (\text{red}, \text{dark}, \text{OW}, \text{No})$

S2 black light overweight no no $P(+|x_1) = P(+|x_2) = P(\text{red}|x_1) \cdot P(\text{OW}|x_1) \cdot P(\text{No}|x_1)$

S3 black dark overweight no yes $= \frac{1}{14} \cdot \frac{3}{5} \cdot \frac{3}{9} \approx 0.01 \approx 0.01$

S4 black light overweight no yes $P(-|x_1) = \frac{5}{14} \cdot \frac{2}{5} \cdot \frac{4}{9} \approx 0.08$

S5 black dark normal yes yes $+ +$

S6 black light normal yes yes $+ +$

S7 black dark overweight yes yes $+ +$

S8 black light overweight yes yes $+ +$

S9 black dark normal yes yes $+ +$

S10 black light overweight yes yes $+ +$

S11 black light normal yes yes $+ +$

S12 black light overweight yes yes $+ +$

S13 black light normal yes yes $+ +$

S14 black light overweight yes yes $+ +$

PPMI(w, c) = $\max(\text{PMI}(w, c), 0)$

independent => not cont

measure informative? -> PMI

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont

PPMI(w, c) + max(PPMI(w, c), 0)

indep. => not cont