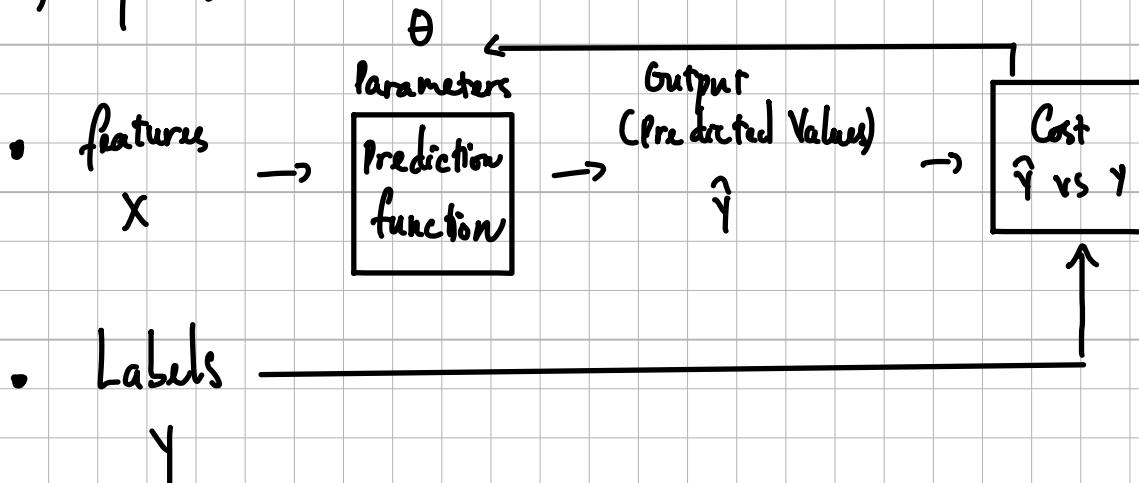


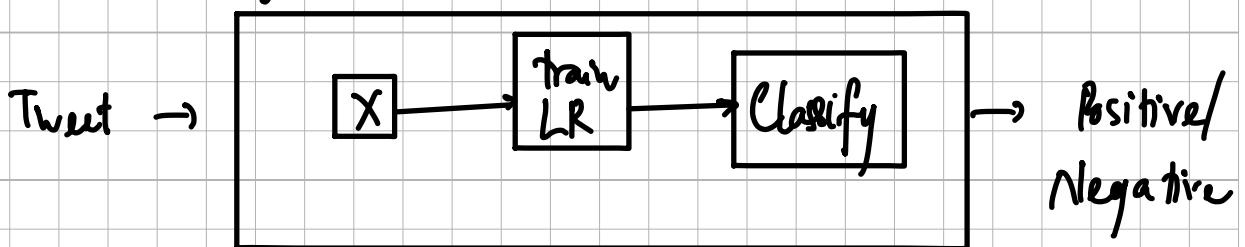
Coursera - NLP

(1) Classification and Vector Spaces

+ Supervised ML



+ Sentiment Analysis



+ Vocabulary

$$V = \text{list of unique words} \quad |V| = n$$

+ Feature Extraction

$$\text{Sparse representations} \rightarrow \theta = [\theta_0 \dots \theta_n]$$

bias
↑

+ Negative and Positive frequencies:

+ Feature Extraction with frequencies:

$$X_m = [1, \sum_w \text{freqs}(w, 1), \sum_w \text{freqs}(w, 0)]$$

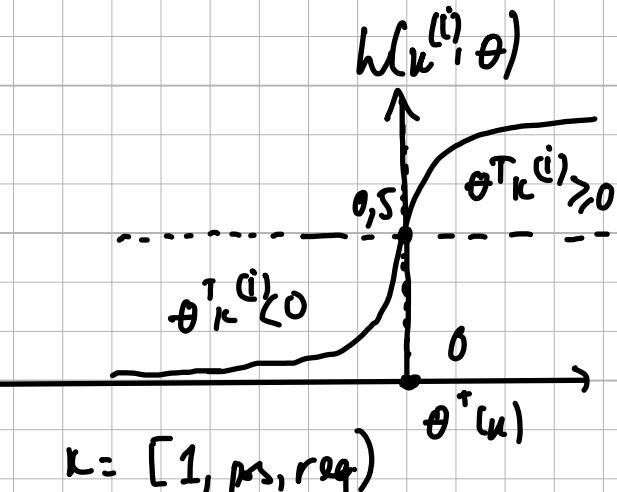
+ Stemming and Stop words

+ Putting it All Together

$$X\text{-shape} = (m, l)$$

$$X = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} \\ \vdots & \vdots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} \end{bmatrix}$$

$$\begin{cases} \theta^T x^{(i)} \geq 0 \\ \theta^T k^{(i)} < 0 \end{cases}$$



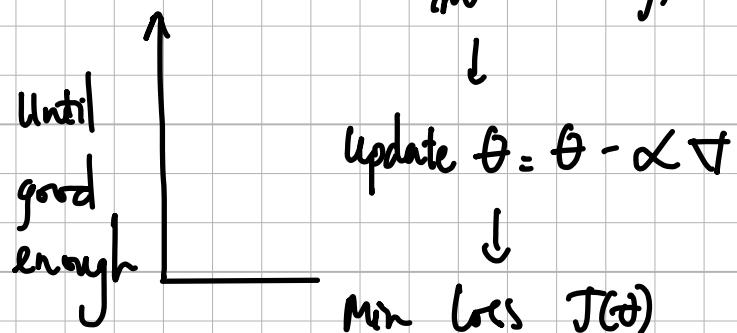
+ Logistic Regression

$$h(x^{(i)}, \theta) = \frac{1}{1 + e^{-\theta^T x^{(i)}}}$$

$$k = [1, pos, neg]$$

+ Training

$$\text{initialize } \theta \rightarrow h = h(X, \theta) \rightarrow \nabla = \frac{1}{m} X^T (h - y)$$



+ Testing

$$\frac{\sum_{i=1}^m (\text{pred } \hat{y}^{(i)}) = y_{\text{val}}^{(i)}}{m}$$

$$\begin{array}{ll} \hat{y}^{(i)} = 0 & \hat{y} = 0 \\ \hat{y}^{(i)} = 1 & \hat{y} = 1 \end{array}$$

+) Cost function Intuition:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log h(x^{(i)}, \theta) + (1-y^{(i)}) \log (1-h(x^{(i)}, \theta))]$$

(2) Naive Bayes

+) Probabilities

→ To calculate a probability of a certain event happening

+) Bayes rule

Conditional Probabilities

$$P(\text{Positive} | \text{"happy"}) = P(\text{"happy} | \text{Positive}) \cdot \frac{P(\text{Positive})}{P(\text{"happy})}$$

$$P(\text{"happy} | \text{Positive}) = 25\%$$

$$P(\text{"happy}) = 15\%$$

$$P(\text{Positive}) = 40\%$$

$$P(X|Y) = P(Y|X) \cdot \frac{P(X)}{P(Y)}$$

$$P(\text{Positive} | \text{"happy"}) = 97\%$$

+) Naive Bayes for Sentiment Analysis

"Events are all independent → rarely the case"

+) Inference condition rule for binary Classification

$$\prod_{i=1}^m \frac{P(w_i | \text{pos})}{P(w_i | \text{neg})},$$

$$\begin{matrix} \gamma_1 & \rightarrow & p_1 \\ \square & \square & \square & \square \end{matrix}$$

+) Laplacean smoothing \downarrow^0 class $\in \{\text{Positive, Negative}\}$

$$\underset{j}{\frac{P(w_i | \text{class})}{\sum}} = \frac{\text{freq}(w_i, \text{class}) + 1}{N_{\text{class}} + V_{\text{class}}}, \rightarrow 0$$

N_{class} + V_{class}
 freq of all words # of unique words in class,

+) log likelihood

≤ 50

positive $\rightarrow \infty$ ratio $(w_i) : \frac{P(w_i | \text{pos})}{P(w_i | \text{neg})},$ Naive Bayes' Reference

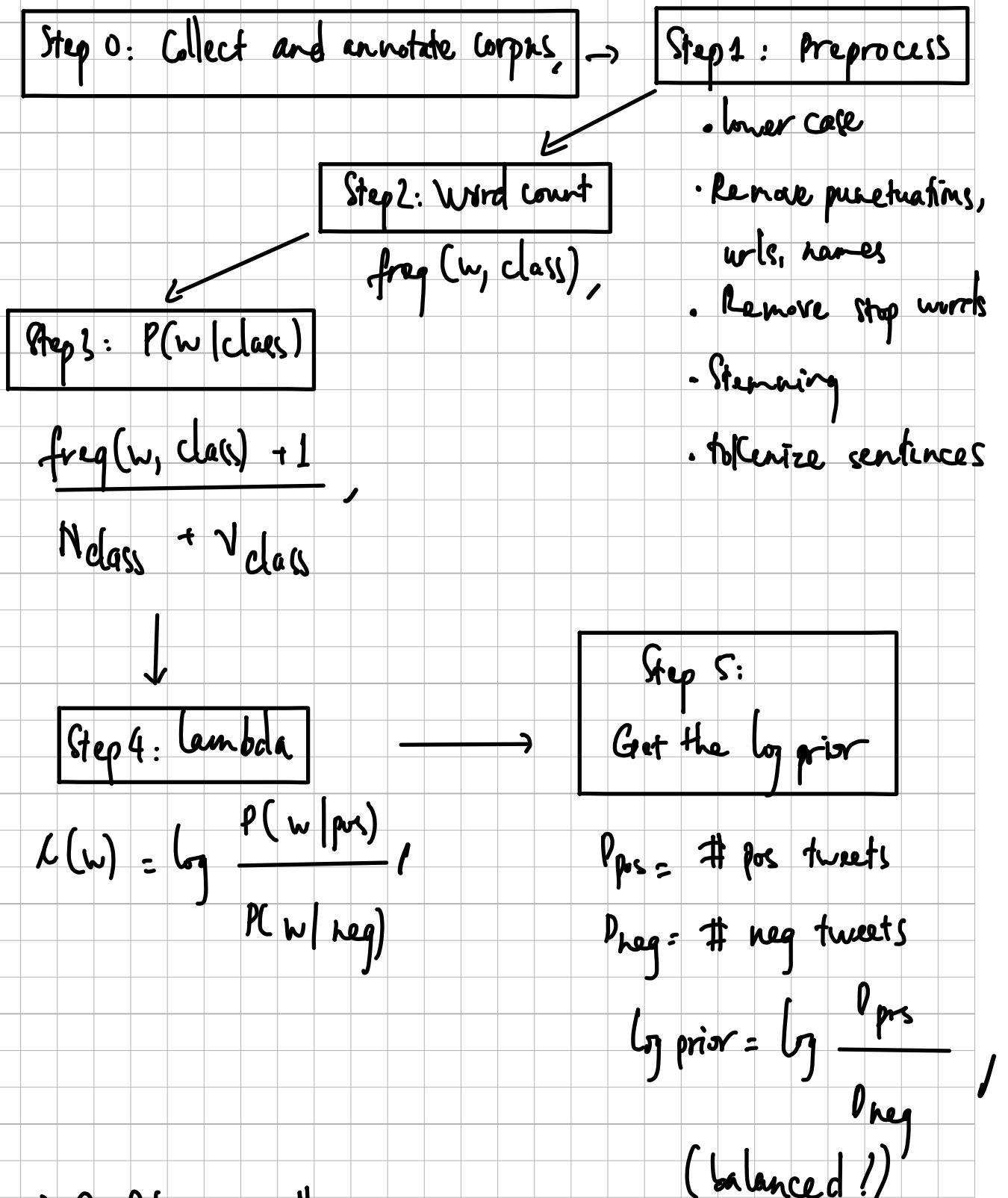
Neutral $\rightarrow 1$
 Negative $\rightarrow 0$ $\approx \frac{\text{freq}(w_i, 1) + 1}{\text{freq}(w_i, 0) + 0} \stackrel{?}{\geq} \log \left(\frac{P(\text{pos}) \cdot \prod_{i=1}^m P(w_i | \text{pos})}{P(\text{neg}) \cdot \prod_{i=1}^m P(w_i | \text{neg})} \right) > 1$

= log prior $\in \{\text{log likelihood}\}, \in [-\infty, \infty]$

Word sentiment $\left\{ \begin{array}{l} \text{ratio}(w) = \frac{P(w | \text{pos})}{P(w | \text{neg})} \\ N(w) = \log \frac{P(w | \text{pos})}{P(w | \text{neg})} \end{array} \right.$

$$\sum_{i=1}^m \log \frac{P(w_i | \text{pos})}{P(w_i | \text{neg})} = \sum_{i=1}^m N(w_i) \quad \geq 0 \rightarrow \text{positive} \\ < 0 \rightarrow \text{negative}$$

→) Training Naive Bayes



+) Confidence ellipse

f) Testing

• X_{val} Y_{val} λ log prior

Score = predict (X_{val} , λ , log prior)

pred = score > 0

$$\frac{1}{m} \sum_{i=1}^m (\text{pred}_i := Y_{\text{val}, i})$$

• Unseen = neutral words

• Application of Naive Bayes

- Sentiment analysis
- Author identification
- Information retrieval
- Word disambiguation
- Simple, fast and robust

+) Naive Bayes Assumptions

- Independence
- Relative frequencies in corpus /

+) Error Analysis

- Removing punctuation and stop words
- Word order

- Adversarial attacks

↳ Sarcasm, Irony and Euphemisms

(3) Vector Space Models



Different meaning vs. Same meaning

Capture dependencies between words,

Information Extraction	Machine Translation	Chatbots
------------------------	---------------------	----------

fundamental concept

"You shall know a word by the company it keeps"

→ neighbouring words → center word

↳ Synonyms and antonyms can be easily interchanged,

+) Word by Word and Word by Doc

C-occurrence matrix

W-by-W: # times occur together within a certain distance k

e.g. single raw like T
data 2 1 1 0

W-by-D: # times occurs within a certain category

e.g. data film	Entertainment 600	Economy 6620	ML 9320
----------------------	----------------------	-----------------	------------

Measure of "similarity": Angle Distance

axis = 0 \rightarrow get for each column

axis = 1 \rightarrow get for each row

+ Euclidean Distance

np, libalg, norm ($v-w$) \Rightarrow problems!

+ Cosine Similarity

/ /

When corpora are different sizes

$\cos = 0 \Rightarrow$ dissimilar $\theta = 90^\circ$

$\cos \rightarrow 1 \Rightarrow$ similar

+ Manipulating Words in Vector Spaces $\rightsquigarrow []$

+ Visualization and PCA $\leftarrow []$

• Motivation

Original space \rightarrow Uncorrelated features \rightarrow Dimensionality reduction

• Visualization to see words' relationships in the vector space.

+ PCA algo: Eigen values, Eigen vectors

Eigenvector: Uncorrelated features for your data

from covariance matrix

Eigenvalue: the amount of information retained by each feature

Mean Normalize
Data

$$x_i = \frac{x_i - \mu_{x_i}}{\sigma_{x_i}}$$



Get Covariance
Matrix

$$\Sigma$$



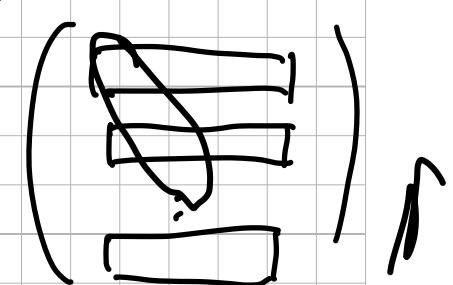
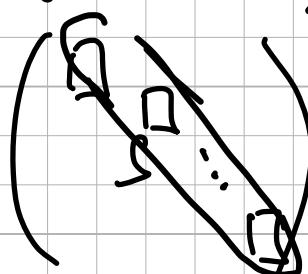
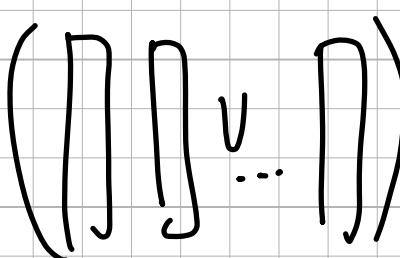
Perform SVD

SVD(Σ)

- Eigen Vectors

Eigenvalues

desc order



Dot product to

Project Data



Percentage of
Retained Variance

$$X' = XU[:, 0:k] \quad \text{or } k-d$$



$$\frac{\sum_{i=0}^k \lambda_i}{\sum_{i=0}^d \lambda_i}$$

Summary :

- Eigen vectors give the direction of uncorrelated features
- Eigenvalues are the variance of the new features
- Dot product gives the projection on uncorrelated features

(4) Machine Translation

Nearest neighbor search using locally sensitive hashing

Machine translation

Document search

+ Transforming word vectors

Align word vectors

$$xR \approx y$$

- → Train on the subset, predict the entire

Solving for R:

initialize R

in loop:

→ Norm 2

$$\text{loss} = \|xR - y\|_F$$

$$g = \frac{d}{dR} \text{loss} \quad \text{gradient}$$

$$R = R - \alpha g \quad \text{update}$$

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}$$

$$\text{np. sqrt}(\text{np. sum}(\text{np. square}(A)))$$

Problem Set (S219)

Gradient

$$\text{loss} = \|X\beta - Y\|_F^2$$

$$g = \frac{d}{d\beta_i} \text{loss} = \frac{2}{n} (X^T (X\beta - Y)) +$$

- +) K nearest neighbors : absolute values, closest
- +) Hash table and hash functions

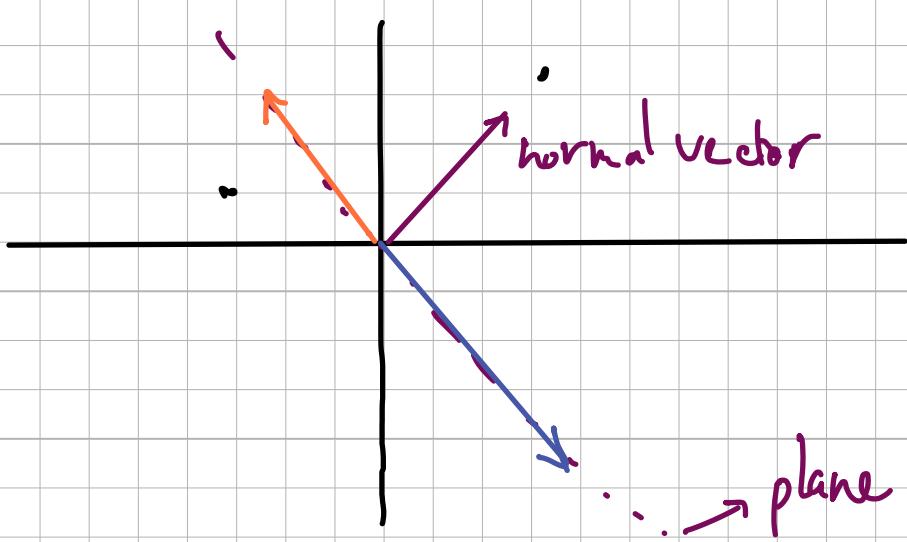
hash function (vector) \rightarrow hash value

- +) locality sensitive hashing
based on where located in vector space

n-buckets

{i: []}

a line = a plane, parallel on the line



def side - of - plane (P, \vec{v}) \rightarrow get relative pos to the
normal vector
plane

+) Multiple planes \rightarrow Dot products \rightarrow hash values

$$\text{sign}_i > 0 \rightarrow h_i = 1$$

$$\text{sign}_i < 0 \rightarrow h_i = 0$$

$$\text{hash} = \sum_i 2^i \cdot h_i \quad G_i: [0..]$$

+) Approximate nearest neighbors

Random planes

mp. random . world (size = (num_planes,
num_dims))

+ Searching documents

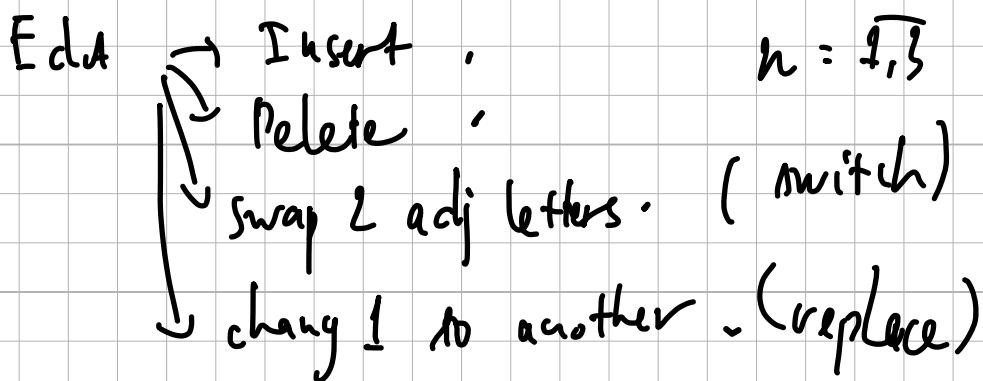
$$\text{doc} = \sum \text{word}$$

(5) Autocorrect and Minimum Edit Distance

1) Autocorrect

- I identify a misspelled word
- find strings in edit distance away
- filter candidates
- calculate word probabilities

+ building the model



Probabilities # times the word appears

$$P(w) = \frac{C(w)}{V} \rightarrow \text{total size of corpus (not unique words)}$$

+ Minimum edit distance | {switch}

• Evaluate similarity between 2 strings

• find the minimum # of edits between 2 strings

Implement spelling corrector, document similarity, machine translation, DNA sequencing and more

$D[i:j] = \text{source}[:i] \rightarrow \text{target}[:j]$,

$D[i:i]$:

$$\min \left\{ \begin{array}{l} D[i-1:j] + \text{del-cost}, \\ D[i, j-1] + \text{ins-cost}, \\ D[i-1, j-1] + \begin{cases} \text{rep-cost}; & \text{if } \text{src}[i] \neq \text{tar}[j] \\ 0; & \text{if } \text{src}[i] == \text{tar}[j] \end{cases} \end{array} \right.$$

• levenshtein distance

• Backtrace

• NP

#	H	Y	E
0			
1			
2			
3			
4			
i			
e			

(c) Part of speech tagging

+) part of speech refers to category of words, lexical terms in the language : adv, v, n, adj...

lexical term	tag
noun	NN
verb	VB
determiner	DT
w-adverb	WRB

the, a
why, where

+) Applications → probs of PoS tags

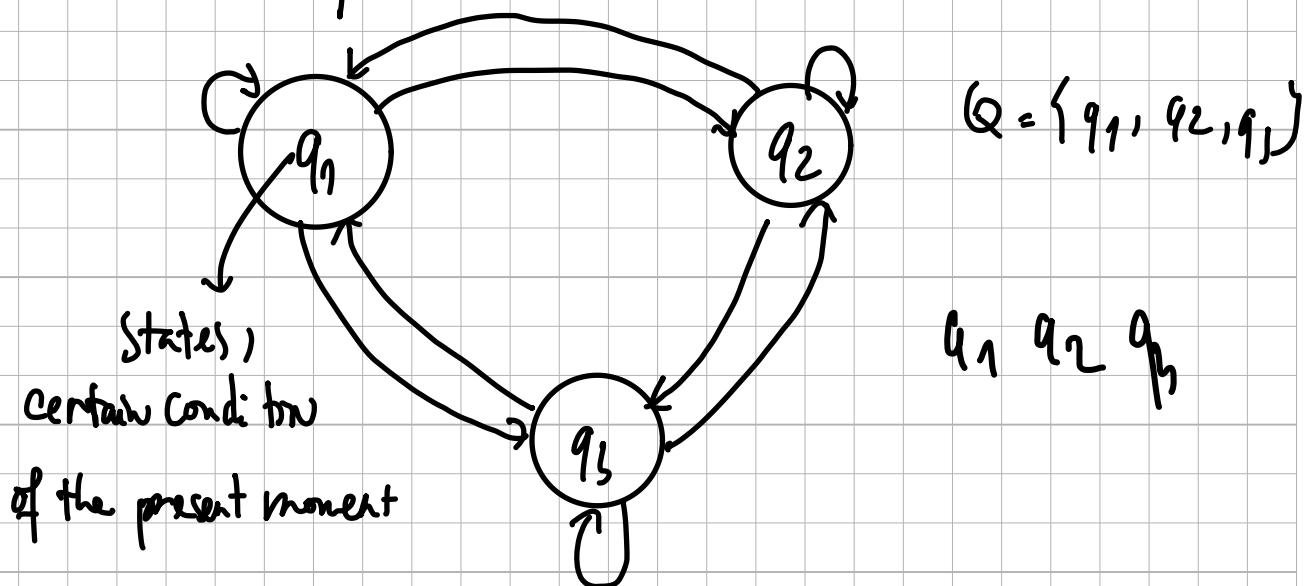
Named entities

Co-reference resolution

Speech recog

+ Markov chains : transition probs, states

Visual representation



+) PoS tags as states 1

Markov property \rightarrow simple, only depends on previous

Transition probabilities (matrix)

current state				# states $(N \times N)$
	NN	V/B	O	
A = <u>$\begin{bmatrix} NN \\ V/B \\ O \end{bmatrix}$</u>	0, 2, 1	0, 1, 2	0, 6	
V/B	0, 4	0, 3	0, 1	
O (other)	0, 2	0, 3	0, 5	
				$\begin{pmatrix} & & \\ & \ddots & \\ \vdots & & \ddots \\ & \ddots & & \ddots \\ - & & & \ddots \end{pmatrix}$
$\sum_{i=1}^N a_{ij} = 1$	(each row)			

? No previous words \rightarrow initial state $(N+1) \times N$

n (initial) 0, 4 0, 1 0, 5

+) Hidden Markov Models \rightarrow decode hidden states

emission prob. $N \times V$ 

\hookrightarrow probability to go from POS tag to a specific word

+) Calculate Probs

• Transition

1. Count occurrences of tag pairs. $\rightarrow ?$:

$$C(t_{i-1}, t_i)$$

$$2. P(t_i | t_{i-1}) = \frac{C(t_{i-1}, t_i)}{\sum_{j=1}^N C(t_{i-1}, t_j)}$$

Preparation of the corpus

$\langle s \rangle$

lowercase - insensitive punctuation

$\langle s \rangle$

- - -

• Smoothing

small value

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

$$P\left(\frac{t_i}{c_j} \mid t_{\text{raw}}\right) = \frac{C(t_{\text{raw}}, t_{c_j}) + 0,001}{\sum_{c_l=1}^N C(t_{\text{raw}}, t_{c_l}) + N + 0,001}$$

+) Emission

You: 2

: 3 → Emission prob : You = 2/3

You eat
the oatmeal
You eat

in a ...

b = NN
VB

j

$$P(w_i \mid t_i) = \frac{C(t_i, w_i) + \epsilon}{C(t_i) + N + \epsilon}$$

$$\sum_{j=1}^V C(t_i, w_j)$$

+) Viterbi Algorithm - graph algo

1. Initialization step

2. forward pass

3. backward pass

prob max each path, keep track path end up taking
 C D

Initialization

(1) Cprob of every word belongs to certain pos

$$w_i: c_{i,1} = \pi_i * b_i^{\text{col}} \text{, index}(w_i)$$

$$= a_{1,i} * b_i^{\text{col}} \text{, index}(w_i)$$

(2) D: store different states you are going through when finding most likely pos tags for the given sequence words $w_1 \dots w_k$

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

+ Forward step

$$c_{ij} = \max_k c_{k,j-1} + a_{k,i} * b_i^{\text{col}} \text{, index}(w_j)$$

$$d_{ij} = \arg \max_k c_{k,j-1} + a_{k,i} * b_i^{\text{col}} \text{, index}(w_j)$$

+ Backward pass

- Index start with 0!

- $\log(c_{i,j}) = \max_k \log(c_{k,j-1}) + \log(a_{k,i}) \log(b_i^{\text{col}} \text{, index}(w_j))$

$$S = \arg \max_i c_{i,L} = \text{index in } D \text{ (same path)}$$

↳ and then use D

(7) Auto Complete

physical talk
↑ or sign

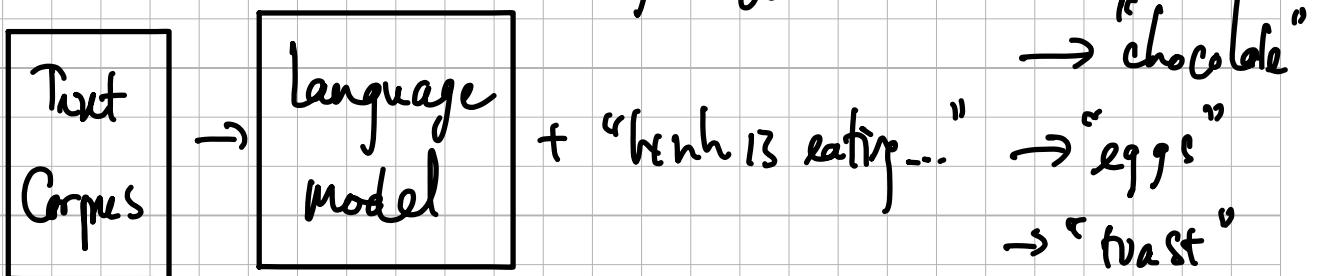
+) N-grams:

speech recog, spelling correction, augmentative communication,..

- Creating language model (LM) from text corpus & :

-) Estimate probability of word sequences
-) Estimate probability of a word following a sequence of words

- Apply this concept to autocomplete a sentence with most likely suggestions



- Sentence auto-complete
- Process a text corpus to N-gram LM
 - Handle out of vocab words
 - Implement smoothing for previously unseen N-grams
 - LM evaluation

+) N-grams and probabilities (conditional prob)

↓
Set of words, punctuation treated like words

Corpus: I am happy because I am learning

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

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

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

• Sequence notation:

$w_1 \dots w_m$

$w_1^m = w_1 w_2 \dots w_m$

$w_1^3 = w_1 w_2 w_3$

$w_{m-2}^m = w_{m-2} w_{m-1} w_m$

Probability ↗

Unigram: $P(w) = \frac{C(w)}{m}$

Bigram: $P(y|x) = \frac{C(xy)}{\sum_w C(xw)} = \frac{C(xy)}{C(x)}$

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

N-gram

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

→ Sequence Probabilities → A appears just before B

$$P(A, B) = P(A) \cdot P(B | A)$$

$$P(A, B, C, D) = P(A) \cdot P(B | A) \cdot P(C | A, B) \cdot P(D | A, B, C),$$

- Probability of a sequence:

$$P(\text{the teacher drinks tea}) \stackrel{\text{chain rule}}{=} \dots$$

- Sentence not in corpus

→ Approximation:

◦ The Markov assumption indicates that only the last word matters

- Bigram:

$$P(w_n | w_{n-1}) \approx P(w_n | w_{n-1}) \quad \begin{array}{l} \text{last-start} \\ \rightarrow \\ +1 \end{array} = N-1$$

Entire sentence

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

$$\approx P(w_1) \cdot P(w_2 | w_1) \dots P(w_n | w_{n-1})$$

- N-gram

$$P(w_n | w_{n-1}^{n-1}) \approx P(w_n | w_{n-N+1}^{n-1})$$

+) Start and ending sentence

$\langle s \rangle$

$\langle /s \rangle$

- N-gram model : add N-1 start token $\langle s \rangle$ at the beginning of each sentence just add one $\langle /s \rangle$

+ Calculate count matrix \rightarrow prob matrix $\rightarrow LM$

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

→ Count matrix

Rows: unique corpus (N-1)-grams

Columns: unique corpus words

$C(w_{n-N+1}^{n-1})$

sum(row)

• Bigram

count matrix \rightarrow Sliding window 2-size

→ Probability matrix : cell (row(sum))

• Language model

- prob matrix $\Rightarrow LM$

+ Sentence prob

+ next word prediction

- Log prod

$$\log(a \cdot b) = \log(a) + \log(b)$$

Algorithm:

1. Choose sentence start (higher value, randomly)
2. Choose next bigram starting with previous word
3. Continue until </s> is picked

+ LM evaluation

- Train / validation / test split

smaller corpora

large corpora (typical for text)

88% train

1% validation

1% test

calculate \leftarrow 88% train

loss
finetune \leftarrow 10% validate

hyperparam

evaluate \leftarrow 1% test

test data - split method

• Continuous text

• Random short sequences

+ Perplexity metric

$$PP(W) = P(S_1, S_2, \dots, S_m)$$

↓
test set

↓
sentence ending with </s>

-1/m

all words in W including
</s> but not <s>

→ measure complexity
of sample of texts

containing m sentences s

- Smaller perplexity = better model

good model: $PP(W) : 60 \rightarrow 20$ (or lower)

- Character level models $PP <$ word-based models PP

- Perplexity for hyphen models

$$PP(W) = \frac{1}{\prod_{i=1}^m \prod_{j=1}^{|S_i|} p(w_j^{(i)} | w_{j-1}^{(i)})}$$

$w_j^{(i)} \rightarrow j^{\text{th}} \text{ word in } i^{\text{th}} \text{ sentence}$

→ Concatenate all sentences in W

$$PP(W) = \frac{1}{\prod_{i=1}^m p(w_i | w_{i-1})}$$

$w_i \rightarrow i^{\text{th}} \text{ word in test set}$

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

- + Out of vocab words

Speech recog
or
question
answering

- Closed vs. Open vocabularies
- Unknown word = out of vocab word (OOV)
- special tag <UNK> in corpus and in input
- flow to create vocab V
 - Max word frequency f
 - Max $|V|$, including words by frequency
 - Use <UNK> sparingly
 - Perplexity - only compare $\log p$ with the same V

+) Smoothing

- Add-one smoothing (Laplacian smoothing)

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

- Add-K-smoothing

$$\frac{\dots + K}{K+V}$$

- Advanced smoothing

- Leeser-Ney smoothing
- Good-Turing Smoothing

+) Backoff

- if N-gram missing \Rightarrow Use (N-1)-gram ...

- 1st discounting e.g. Katz backoff

- "stupid" backoff : $(N-1) * c$ ($c=0, 4$)

e.g. PC chocolate (John drinks)



$0,4 * P(\text{chocolate} | \text{drinks})$

+) Interpolation

$$\hat{P}(\text{chocolate} | \text{John drinks}) = 0,7 * P(\text{chocolate} | \text{sohachiku}) + 0,2 * P(\text{chocolate} | \text{drinky}) + 0,1 * P(\text{chocolate} | \text{other})$$

$\sum_i \lambda_i = 1 \rightarrow$ optimize lambda using validation set

(8) Word Embeddings

a) Applications

- Semantic analogies and similarity
- Sentiment analysis
- Classification of customer feedback
- Machine translation
- Information extraction
- Question answering

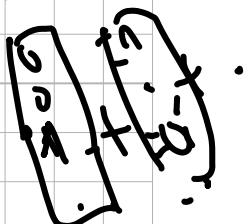


→ Continuous Bag of words or CBOW

- try to predict the particular word given a few current text words
- the weights used to predict \rightarrow word vectors

$\begin{matrix} \text{weight} & \text{weight} \\ \text{for words} & \text{for words} \\ \text{predict} & \text{context} \end{matrix}$

\rightarrow manipulate 2 weights



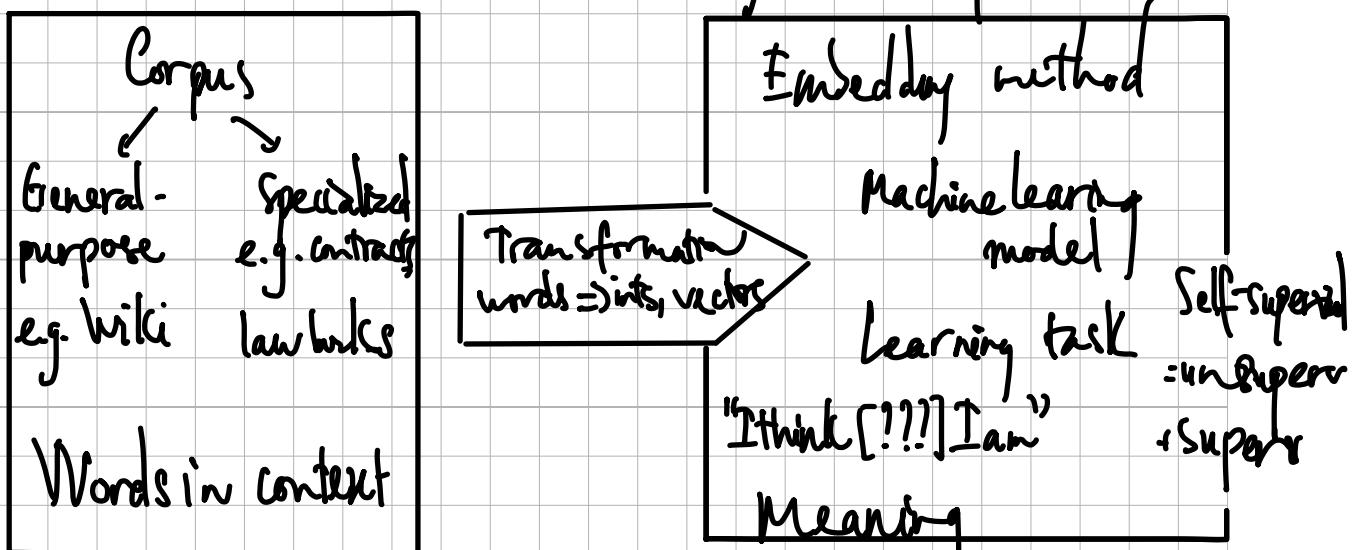
+1 Basic Word representations

- Integers \rightarrow cons: implied ordering not good
- One-hot vectors $\xrightarrow{\quad}$ pros: simple, no implied, cons: huge, encode no meaning
- Word embeddings

- +) Word embeddings
 - low dimension
 - Embeds meaning

e.g. semantic distance
e.g. analogies

- +) How to create Word Embeddings



? What will ultimately define the meaning of the individual words? the specific of task
 => the meaning of words, as carried by the word embeddings, depends on the embedding approach

- +) Word Embedding methods
 - Word2vec (Google, 2013)
 - Continuous Bag-of-Words (CBOW)
 - Continuous skip-gram / Skip-gram with negative sampling (SGNS)

- Global Vectors (GloVe) (Stanford, 2014)
 - fastText (Facebook, 2016)
 - supports out-of-vocabulary (OOV) words
 - Advanced word embedding methods
 - Deep learning, contextual embeddings
 - Tunable pre-trained models available
 - BERT (Google, 2018)
 - ELMo (Allen Institute for AI, 2018)
 - GPT-2 (OpenAI, 2018)
 - + Continuous Bag-of-Words Model
 - Center word prediction: rationale
 - e.g. the little ? is barking
 ↓
 dog puppy howl terrier ...
 - CBOW in a nutshell
- | INPUT | PROJECTION | OUTPUT |
|----------|------------|--------|
| $w(t-2)$ | | |
| $w(t-1)$ | → | $w(t)$ |
| $w(t+1)$ | → | |
| $w(t+2)$ | → | |
- sum → $w(t)$
- context half size of k
- $\frac{1}{k-left} + \frac{1}{k-right}$

objective: to predict a missing word based on surrounding words

+1) Cleaning & Tokenization

- Letter case "the" = "the" = "THE" → lowercase / uppercase

- Punctuation , ! . ? → . " () " → φ ... ! ? ! → .

- Numbers 1 2 ... → φ 1.49159 98260 → ar is / (NUMBER)

- Special characters π \$ € § π ** → φ

- Special words ☺ #nlp → : happy: #nlp

+2) Sliding Window of words

i: c → list of word

while i < len(words) - c: → context size

center = word [i]

context = words[(i-c): i] + words [(i+1): (i+c+1)]

yield context, center # π 14

i += 1

+3) Transform words into Vectors

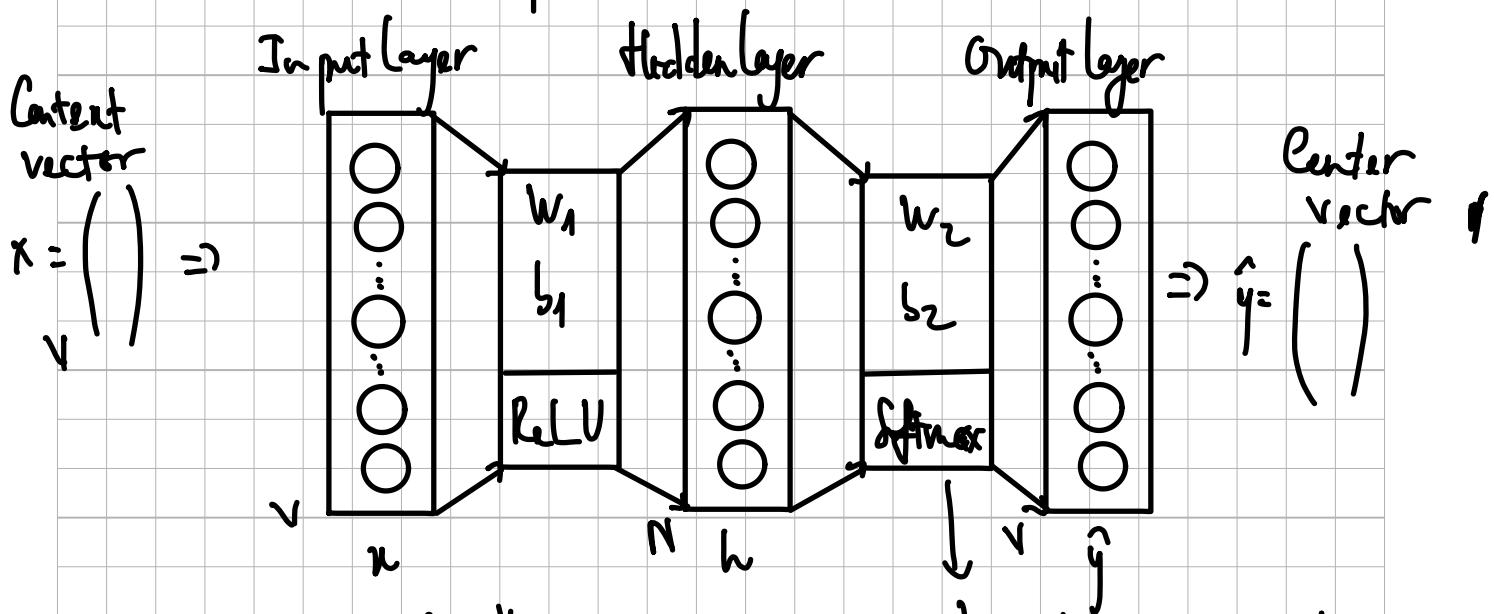
- Started with one-hot vectors for the content words

- transform into a single one = taking average

Content words Content vec Center word Center vec

I am because I [0, 25, 0, 25, 0, 5, 0] happy [0, 0, 1, 0, 0]

f) Architecture of the CROW model



x : average of all content vectors distribution over V
 pick output = argmax of output

+1 Dimensions

$$\begin{array}{ccccccc}
 & w_1 & N \times V & & w_2 & V \times N & \\
 x & \rightarrow & b_1 & N \times 1 & \rightarrow & h & N \times 1 \\
 V \times 1 & & & & & & \\
 & \text{ReLU} & & & \text{softmax} & & \\
 z_1 & = & w_1 x + b_1 & N \times 1 & z_2 & = & w_2 h + b_2 & V \times 1 \\
 & & & & & & & \\
 h & = & \text{ReLU}(z_1) & N \times 1 & \hat{q} & = & \text{softmax}(z_2) & V \times 1
 \end{array}$$

- np.mean(content_words_vectors, axis=0)

- Batch processing

$$\hat{Y} = \left(\left(\hat{y}^{(1)} \right), \dots, \left(\hat{y}^{(m)} \right) \right)$$

Content words vectors

$$\begin{array}{c}
 x = \left(x^{(1)} \right), \dots, \left(x^{(m)} \right) \\
 m \rightarrow \text{batchsize}
 \end{array}$$

$$x : V \times m$$

$$w_1 : N \times V$$

$$b_1 : N \times m$$

$$h = N \times m$$

$$z_2 = N \times m$$

$$w_2 : V \times N$$

$$b_2 : V \times m$$

$$z_2 = V \times m$$

$$\hat{Y} = V \times m$$

$V = 1000$

$m = 16$

$N = 400$

batches of 16
 $W_1: N \times V$ l_1, l_2 broadcast
 $\downarrow \quad \downarrow$
 $N \times 1 \quad V \times 1$

+) Rectified linear unit (ReLU)

$$\text{ReLU}(v_{ii}) = \max(v_{ii}, 0) -$$

Soft Max

$$\begin{matrix} z \\ \vdots \\ z_i \\ \vdots \\ z_N \end{matrix} \xrightarrow{\text{softmax}} \begin{matrix} \hat{y}_1 \\ \vdots \\ \hat{y}_i \\ \vdots \\ \hat{y}_N \end{matrix}$$

Probabilities of being center word

$\hat{y}_i = \frac{e^{z_i}}{\sum_{j=1}^N e^{z_j}}$

+) Training a CNN model

Cat function

↳ cross-entropy loss \rightarrow classification model

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

$$\text{wrong : } J = - \log \hat{y}_k$$

incorrect predictions:

penalty

correct predictions:

reward

+ forward propagation

log \rightarrow single example

Cost \rightarrow a batch of examples = mean of losses.

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

+ Backpropagation and Gradient Descent

$$\frac{\partial J_{\text{batch}}}{\partial W_1} = \frac{1}{m} (W_2^T (\hat{Y} - Y) \cdot \underbrace{\text{step}(Z_1)}_{\text{ReLU}(W_2^T (\hat{Y} - Y))} X^T)$$

$$\frac{\partial J_{\text{batch}}}{\partial W_2} = \frac{1}{m} (\hat{Y} - Y) \underbrace{H^T}_{\text{ReLU}(W_2^T (\hat{Y} - Y))}$$

$$\frac{\partial J_{\text{batch}}}{\partial b_1} = \frac{1}{m} (W_2^T (\hat{Y} - Y) \cdot \underbrace{\frac{1}{m} \text{step}(Z_1)}_{\begin{array}{l} \uparrow \\ \text{np.sum}(a, axis=1, keepdims=True) \end{array}})$$

$$\frac{\partial J_{\text{batch}}}{\partial b_2} = \frac{1}{m} (\hat{Y} - Y) \underbrace{1_m^T}_{\frac{1}{m}}$$

$$\text{Minimize: } W := W - \lambda dW$$

$$b := b - \lambda db$$

f) Extracting Word Embedding Vectors

option 1: $W_1: N \times V$ $x = (V, 1)$

→ each col = w each word

Option 2: $W_2: V \times N$ $x = (1, V)$

→ each row = w each word

option 3

$W_3 = 0.5 (W_1 + W_2^\top) = N \times V$

→ each col = w each word

f) Evaluating Word Embeddings

Intrinsic Evaluation

Continuous skip-grams model vs. Continuous bag-of-words model

→ predict contexts given center | predict target given surroundings
allows you to test relationships between words

| semantic analogies, syntactic analogies.

? harder to track: Ambiguous cases

→ • Analogies

• Clustering

• Visualization

→ Word embeddings were created by Continuous skip-grams, not Contd bag-of-words

Extrinsic Evaluation

Test word embeddings on external task

e.g named entity recognition, parts-of-speech tagging

+ Evaluates actual usefulness of embeddings

- True - encouraging

- More difficult to trouble shoot

library: Keras, Pytorch

(9) Recurrent Neural Networks for Language Modeling

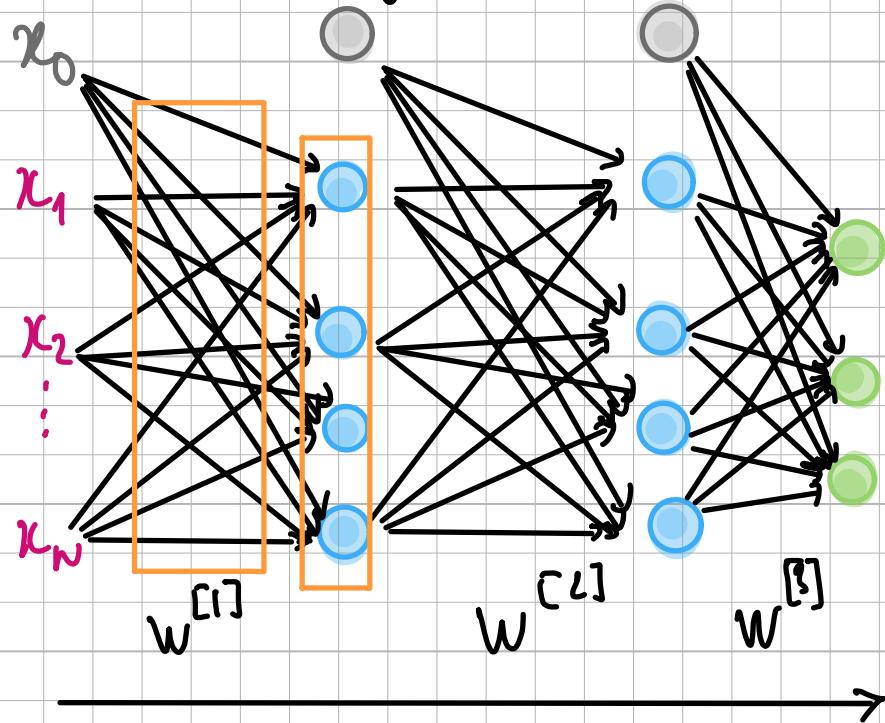
†) Introduction

Deep neural networks can learn more abstract information and do not require any feature engineering.

†) Neural Networks for Sentiment Analysis

- Neural networks and forward propagation

Forward propagation



$a^{[i]}$ Activations in layer i

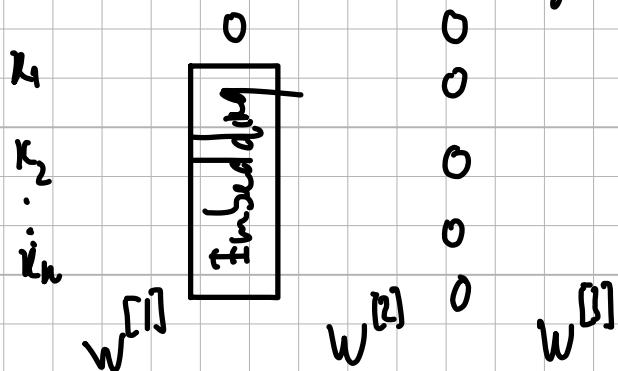
$$a^{[0]} = X$$

$$z^{[i]} = W^{[i]} \cdot a^{[i-1]}$$

$$a^{[i]} = g^{[i]}(z^{[i]})$$

↑ activation func

- Sentiment analysis



0 - -> Positive (1)
0 - -> Negative (0)

f) Initial Implementation
+ padding

+ Dense layers and ReLU

→ Using tensorflow

`tf.keras.layers.Dense(32)`

`tf.keras.layers.ReLU()`

→ Dense layer with 32 units,

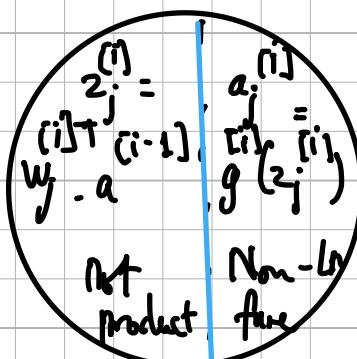
ReLU layer

Dense layer

$$z^{[i]} = \frac{W^{[i]} a^{[i-1]}}{\text{trainable}}$$

parameters

`= tf.keras.layers.Dense(32)`
`activation='relu'`)



$$g(z^{[i]}) = \max(0, z^{[i]})$$

+ Embedding and Mean layers

• Embedding layer

Trainable weights : Vocabulary

x
Embedding

• Mean layer take the average of the embeddings.

Same size as # N (features)

Does not have any trainable param

→ N grams vs. Sequence Models

+) traditional language models

• limitations:

- Large N-grams needed to capture dependencies between distant words

- Need a lot of space and RAM

+) Recurrent Neural Networks



- RNNs model relationships among distant words

- In RNNs a lot of computations share parameters



+) Applications of RNNs

One to One : given scores → predict writer

One to Many < inputs a sentence and generates a sentence
inputs a topic and generates a conversation about
that topic

Many to one : give a tweet → predict sentiment

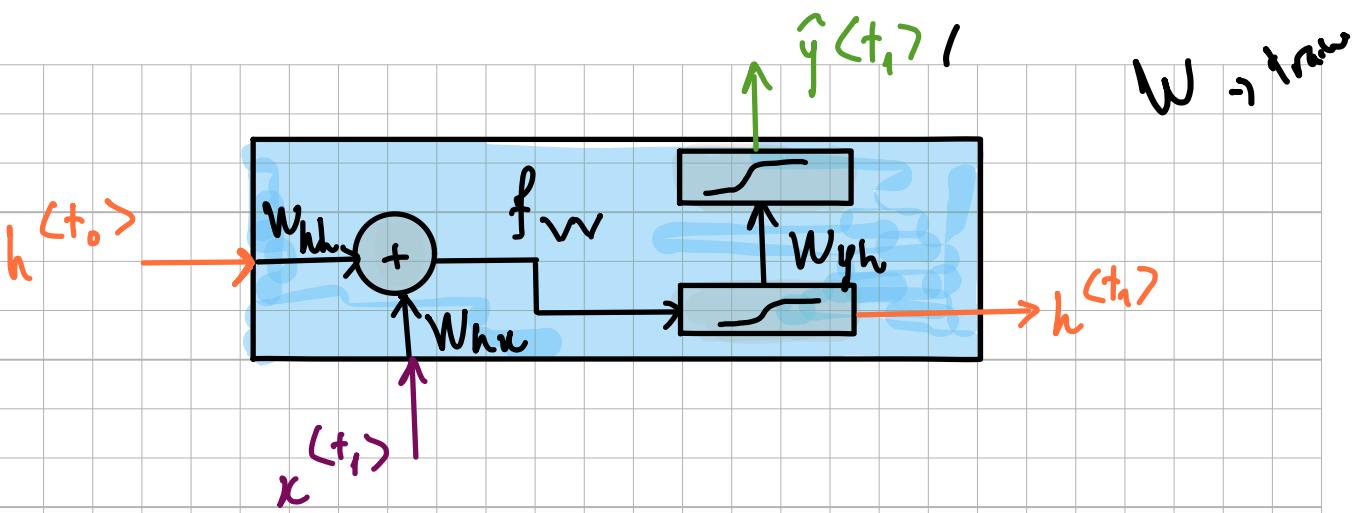
Many to many : machine translation

+) Math in simple RNNs

- Hidden states propagate information through time

- Basic recurrent units have two inputs at each time

$$h^{(t-1)}, h^{(t)}$$



Note that:

$$h^{(t)} = g(W_h[h^{(t-1)}, x^{(t)}] + b_h)$$

is the same as multiplying W_{hh} by h and W_{hx} by x

$$h^{(t)} = g(W_{hh} \overline{h^{(t-1)}} + W_{hx} \overline{x^{(t)}} + b_h)$$

$$\hat{y}^{(t)} = g(W_{yh} \overline{h^{(t)}} + b_y)$$

→ Training W_{hh} , W_{hx} , W_{yh} , b_h , b_y

$$[h \times h, h \times x] \rightarrow h \times (h+x)$$

$$W_h = [W_{hh} | W_{hx}] \text{ horizontal concatenation}$$

→ Option 1: np.concatenate((W_{hh} , W_{hx}), axis=1)

→ Option 2: np.vstack((W_{hh} , W_{hx}))

$$[h^{(t-1)}, x^{(t)}] = \begin{bmatrix} h^{(t-1)} \\ x^{(t)} \end{bmatrix}$$

→ Option 1: up.concatenate(($h_{t_{prev}}$, x_t))

→ Option 2: np.vstack(($h_{t_{prev}}$, x_t))

• Formula 1:

$$\text{np.matmul}(\text{np.vstack}((\mathbf{w}_{hh}, \mathbf{w}_{hw})), \text{np.vstack}((\mathbf{h}_{t-1}, \mathbf{x}_t)))$$

• Formula 2:

$$\text{np.matmul}(\mathbf{w}_{hh}, \mathbf{h}_{t-1}) + \text{np.matmul}(\mathbf{w}_{hx}, \mathbf{x}_t)$$

$\rightarrow \text{np.allclose}(\text{formula_1}, \text{formula_2}) : \text{True}$

+) Cost function for RNNs

for RNNs the loss function is just an average through time

$$J = -\frac{1}{T} \sum_{t=1}^T \sum_{j=1}^k y_j^{(t)} \log \hat{y}_j^{(t)},$$

+) Implementation note:

Scan() function in tensorflow (GPU, compute in parallel)

scan(f_n , x , $h_{(0)}$, initializer=None, ...)

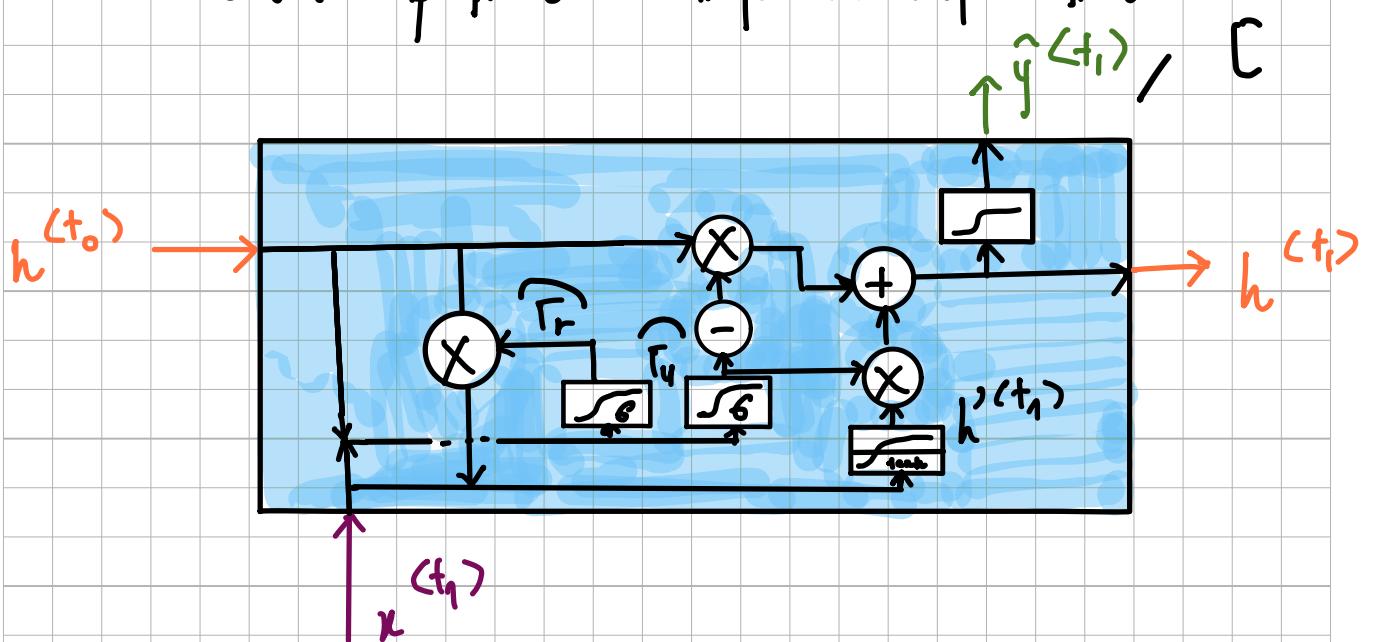
$$f_n \downarrow \quad x^{(1)} \dots x^{(T)} \quad h^{(0)}$$

return y_s , h_T value

+) Gated Recurrent Units (GRU)

Relevance and update gates to remember important
prior information

- GRUs "decide" how to update the hidden state
- GRUs help preserve important information



$$T_u = \sigma(W_u [h^{(t_0)}, x^{(t_1)}] + b_u) \rightarrow \text{update gate}$$

$$T_f = \sigma(W_f [h^{(t_0)}, x^{(t_1)}] + b_f) \rightarrow \text{forget relevance score}$$

$$h^{(t_1)} = \tanh(W_h [T_f * h^{(t_0)}, x^{(t_1)}] + b_h)$$

$$h^{(t_1)} = (1 - T_u) * h^{(t_0)} + T_u * h^{(t_1)}$$

$$\hat{y}^{(t_1)} = g(W_y h^{(t_1)} + b_y)$$

$$c = \text{np.tanh}(c)$$

return_sequences = True

batch_size

sequence_length → not fixed, "None", expect any
 word-vector-length → fixed in model

model_.GRU.build([None, None, word-vector-length])
 .summary()

+ Deep and Bi-directional RNNs ~~→~~

- BiRNN
acyclic graph (không chu trình)

↳ Information flows from the past and from the future independently

$$\hat{y}^{(t)} = g(W_y [\vec{h}^{(t)}, \vec{h}^{(t)}] + b_y)$$

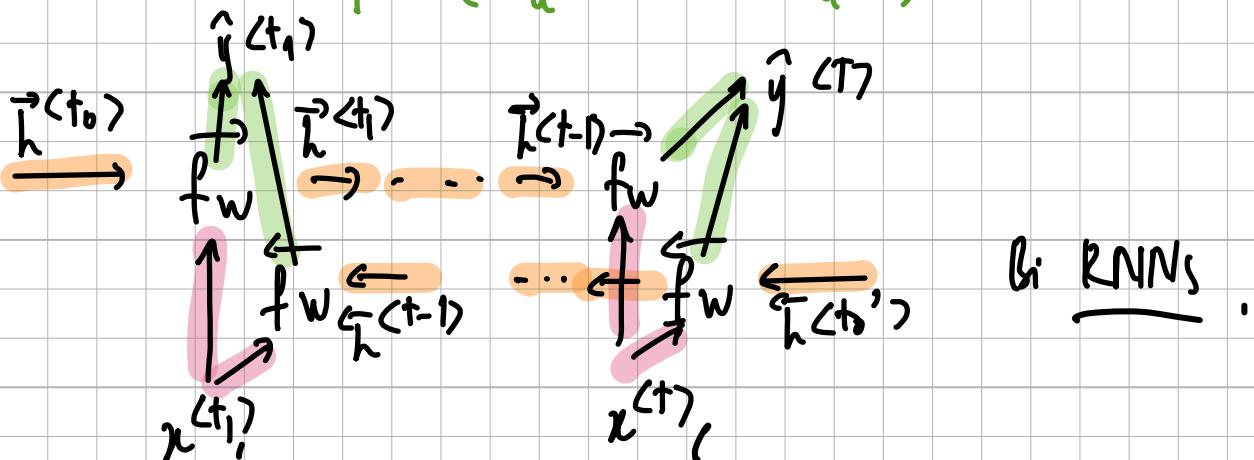
• Deep RNN

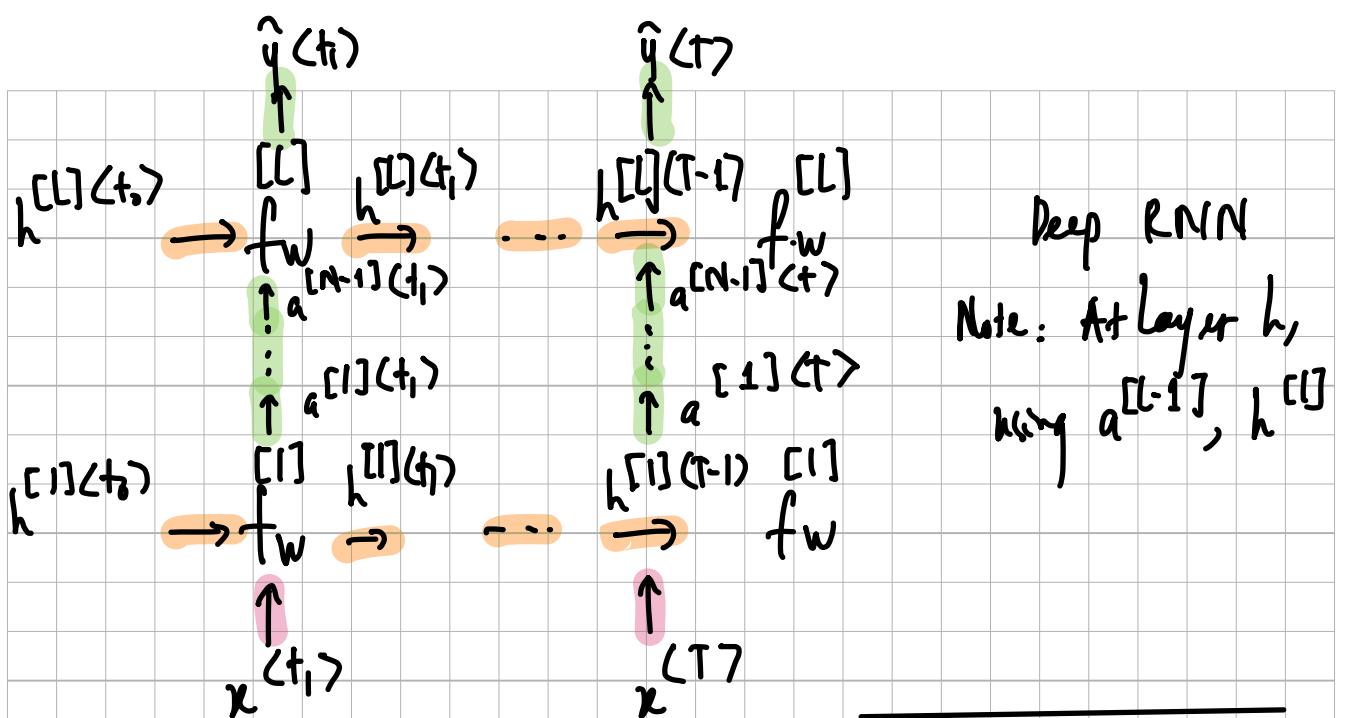
1. Get hidden states for current layer →

$$h^{[l](t)} = f^{[l]}(W_h^{[l]} [h^{[l-1](t-1)}, a^{[l-1](t)}] + b_h^{[l]})$$

2. Pass the activations to the next layer ↑

$$a^{[l](t)} = f^{[l]}(W_a^{[l]} h^{[l](t)} + b_a^{[l]})$$





Deep RNN

Note: At layer \$h\$,
using \$a^{[l-1]}, h^{[l]}

†) Perplexity

$$\log p(W) = -\frac{1}{N} \sum_{i=1}^N \log P(w_i | w_1, \dots, w_{i-1})$$

$$PP(W) = \sqrt{\prod_{i=1}^N \frac{1}{P(w_i | w_1, \dots, w_{i-1})}}$$

(10) LSTMs and Named Entity Recognition

→ RNNs and Vanishing Gradients

Solving for vanishing or exploding gradients

- Identity RNN with ReLU activation

$$\begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix} \quad -1 \rightarrow 0$$

- Gradient clipping $32 \rightarrow 25$

- Skip connections

explode or vanish

$$x \rightarrow \boxed{\quad} \rightarrow \oplus \rightarrow F(x) + x$$

$$\sigma [0, 1]$$

$$\tanh [-1, 1]$$

→ depends on eigen values λ

Max W $\begin{cases} < \lambda \rightarrow \text{vanish} \\ \end{cases}$

$\begin{cases} > \lambda \rightarrow \text{explode} \\ \end{cases}$

Derivative of activation func bounded by $\alpha = \frac{1}{\lambda}$

$$d(\text{sigmoid}) = (\text{sigmoid})(1 - \text{sigmoid}) \in [0, 0.25]$$

+) Intro to LSTMs (long short-term memory unit)

Gates in LSTM

Cell state

$c_{(t_0)}$

\rightarrow



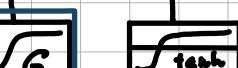
1



2



3



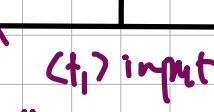
hidden state

$h_{(t_0)}$

\rightarrow

$x_{(t_1)}$ input

\uparrow



$+$



\rightarrow



\downarrow



\uparrow



\uparrow



\uparrow



\uparrow



1. Forget gate (to keep) f

information that is no longer important

2. Input gate (to add) i

information to be stored

3. Output gate (next hidden states will be) o

information to use at current step

Applications

- Next - character prediction (email?)

- Chatbots

- Music composition

- Image Captioning

- Speech recognition

$\hat{y}(t)$ output

$c_{(t_1)}$

$h_{(t_1)}$

+) LSTM architecture

The gate (candidate memory cell):

$$g = \tanh(W_g [h_{t-1}; x_t] + b_g)$$

The cell state:

$$c_t = f \odot c_{t-1} + i \odot g$$

The output of LSTM unit

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

+) Intro & NER

locates and extracts predefined entities from text

Places, organizations, names, time and classes

Types of entities

Geographical

organization

Geopolitical

Time indicator

Artifact

Person

Applications

Search efficiency

Recommendation engines

Customer service

Automatic trading (news...)

f) Training NERs

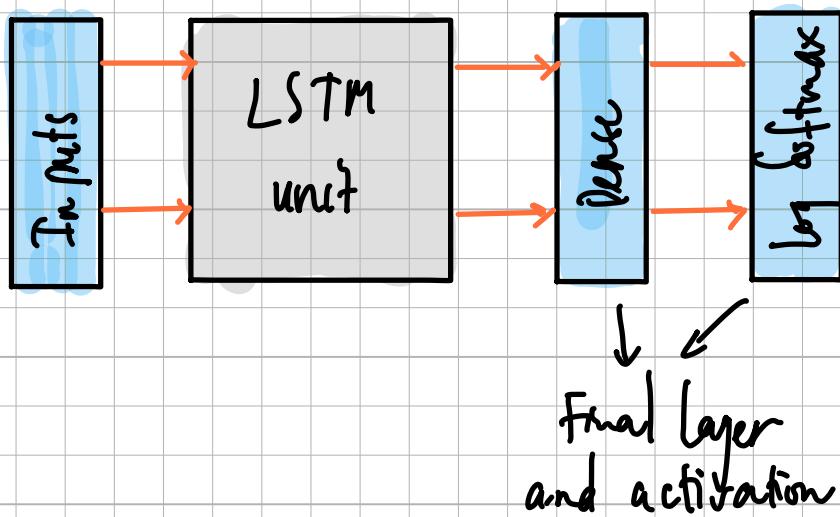
Data Processing

(same length)

- Convert words and entity classes into arrays
- Pad with tokens: Set sequence length to a certain number and use <PAD> token to fill empty spaces
- Create a data generator

Training an NER system

1. Create a tensor for each input and its corresponding number (faster processing)
2. Put them in a batch $\rightarrow 64, 128, \dots$
3. Feed it into an LSTM unit
4. Run the output through a dense layer
5. Predict a log soft max over K classes



g) Computing accuracy

- Pass test set through the model

- Get argmax across the prediction array
- Mask padded tokens
- Compare with true labels

(1) Siamese networks and the triplet loss

I identify similarity between things

similarity

+ Applications:

Handwritten check - \sim

Question duplicates / \sim word.

Queries

+ Architecture

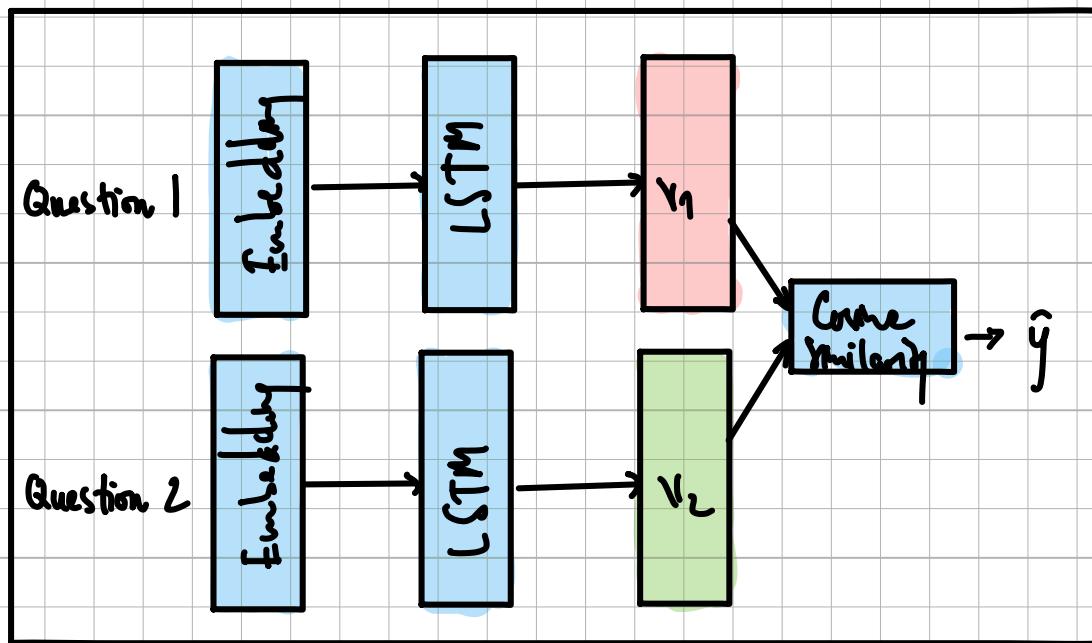
Subnetworks share identical parameters, only need to train one

1. Embedding

2. LSTM (optional)

3. Vectors (output of each sub-network)

4. Cosine Similarity



lab

→ Creating a Siamese model

Import tensorflow as tf

from tensorflow.keras import layers

from tensorflow.keras.models import Model, Sequential

from tensorflow import math

import numpy

Setting random seeds

numpy.random.seed(10)

+) Cost function : Triplet Loss

Anchor , Positive , Negative

$$s(A, P) \approx 1$$

$$s(A, N) \approx -1$$

$$\text{diff} = s(A, N) - s(A, P)$$



$$-\cos(A, P) + \cos(A, N) \leq 0$$

predictor $-\cos(A, P)$ $\cos(A, N)$ sum

- (1)

- 1

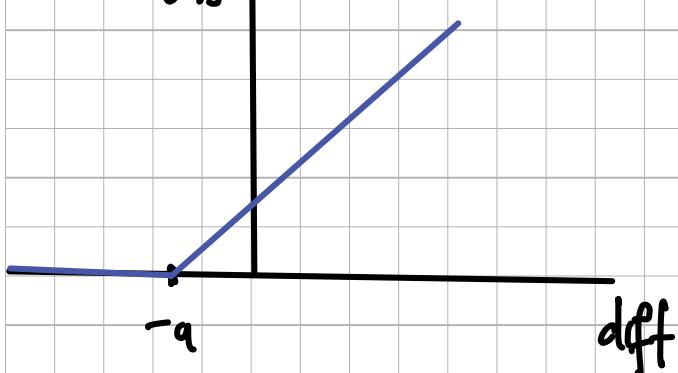
- 2

- (-1)

1

2

+) Triplets loss



with alpha margin

$$L = \begin{cases} 0 & \text{if } \text{diff} + \alpha \leq 0 \\ \text{diff} + \alpha & \text{if } \text{diff} + \alpha > 0 \end{cases}$$

Full cut : $\max(-\cos(A, P) + \cos(A, N) + \alpha, 0)$

- easy neg triplet:

$$\cos(A, N) < \cos(A, P)$$

controls how far
 $\cos(A, P)$ is from $\cos(A, N)$

- Semi-hard neg triplet:

$$\cos(A, N) < \cos(A, P) < \cos(A, N) + \alpha$$

- Hard neg triplet:

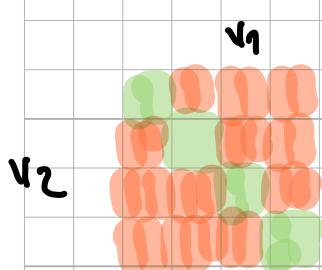
$$\cos(A, P) < \cos(A, N)$$

- +) Computing the cost.

$$s(v_1, v_2)$$

$$L(A, P, N) = \max(\text{diff} + \alpha, 0)$$

$$\text{diff} = s(A, N) - s(A, P)$$



$$J = \sum_{i=1}^m L(A^{(i)}, P^{(i)}, N^{(i)})$$

→ the diag line corresponds to scores of similar sentences (normally they should be positive). The off-diagonals correspond to cos scores between the anchor and the neg samples

2 concepts:

- mean-neg: speeds up training
- closest-neg: helps penalize the cut more

$$L_{\text{original}} = \max(-\cos(A, P) + \cos(A, N + \alpha, 0))$$

$$L_1 = \max(-\cos(A, P) + \text{mean-neg} + \alpha, 0)$$

$$L_2 = \max(-\cos(A, P) + \text{closest-neg} + \alpha, 0)$$

$$L_{\text{full}} = L_1 + L_2$$

+/- One shot learning

Classification vs. One shot learning
classifies input as 1 of K classes

measures similarity between classes

→ no need to retrain, learn a similarity score

$$\cos(sg1, sg2) > T \quad \boxed{\checkmark}$$

$$\cos(sg1, sg2) \leq T \quad \boxed{x}$$

+) Training / Testing

Prepare Batches \rightarrow embedding vector space
 $v_1 = (1, d_model)$

Training

1. Embedding
2. LSTM
3. Vectors
4. Cosine Similarity

Testing

1. Convert two inputs into a array of numbers
2. Feed it into your model
3. Compare v_1, v_2 using cosine similarity
4. test against a threshold T

(12) Neural Machine translation

Adding an attention mechanism to allow the decoder to access all relevant parts of the input sentence regardless of its length

+> Seq2Seq 2014 Google

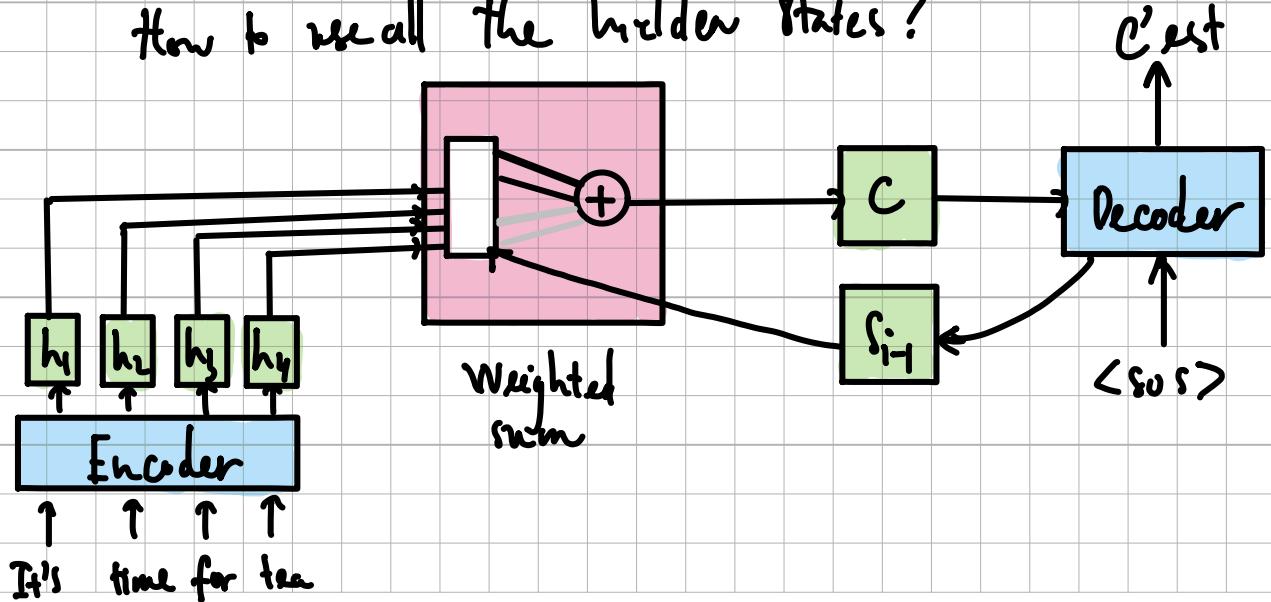
- Maps variable-length sequences to fixed-length memory
- Inputs and outputs can have different lengths
- LSTMs and GRU to avoid vanishing and exploding gradient problems

The information bottleneck : struggles to compress longer sequences

→ New layer : Attention

+> Seq2Seq Model with attention

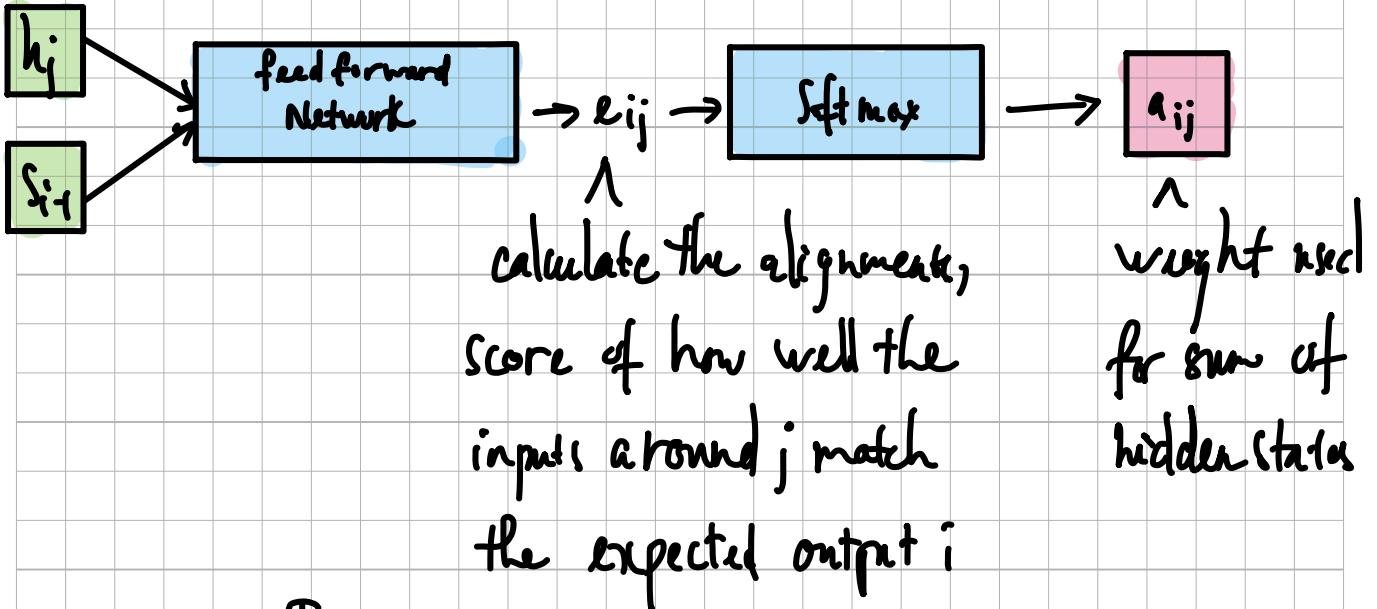
How to recall the hidden states?



The attention layer is more depth

$$e_{ij} = a(s_{i-1}, h_j) \quad \alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^K \exp(e_{ik})}$$

learnable params



$$c_i = \sum_{j=1}^{T_K} \alpha_{ij} \cdot h_j$$

Context Vector is an expected value

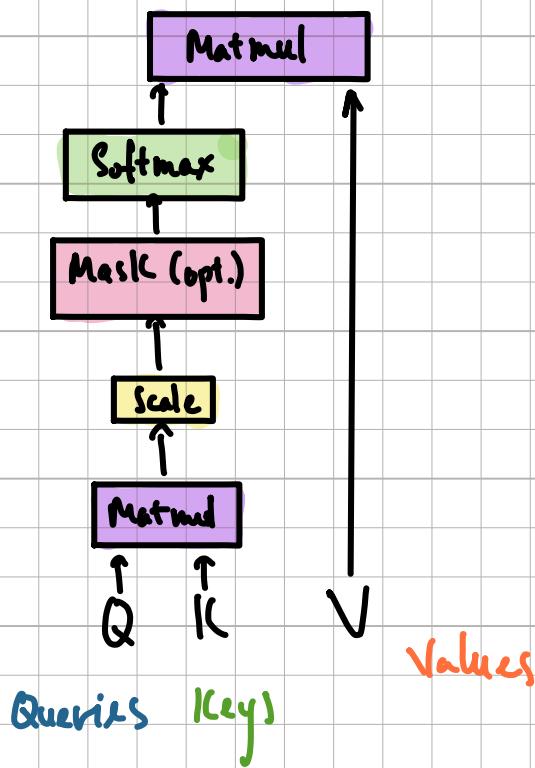
$$\alpha_{i1} \cdot h_1 + \alpha_{i2} \cdot h_2 + \dots + \alpha_{im} \cdot h_m \rightarrow c_i$$

$$e_{ij} = v_a^T \tanh(W_a s_{i-1} + U_a h_j)$$

+) Queries, Keys, Values, and Attention

The similarity between words is called alignment of how well query and keys match

Scaled dot-product attention



weight for weighted sum

$$\text{softmax} \left(\frac{QK^T}{\sqrt{dk}} \right) v$$

Attention

Scale using the root of the key size

→ alignment is learnt in the input embeddings other linear before attention layer

Flexible attention

Work for languages with different grammar structures!
despite different ordering

- +1 Set up for machine translation
 - Use pretrained vector embeddings
 - Otherwise, initially represent words with one-hot vectors
 - Keep track of index mappings with word2ind and ind2word dictionaries
 - Add end of sequence tokens: (EOS)
 - Pad the token vectors with zeros

a) Teacher forcing

Errors from early steps propagate

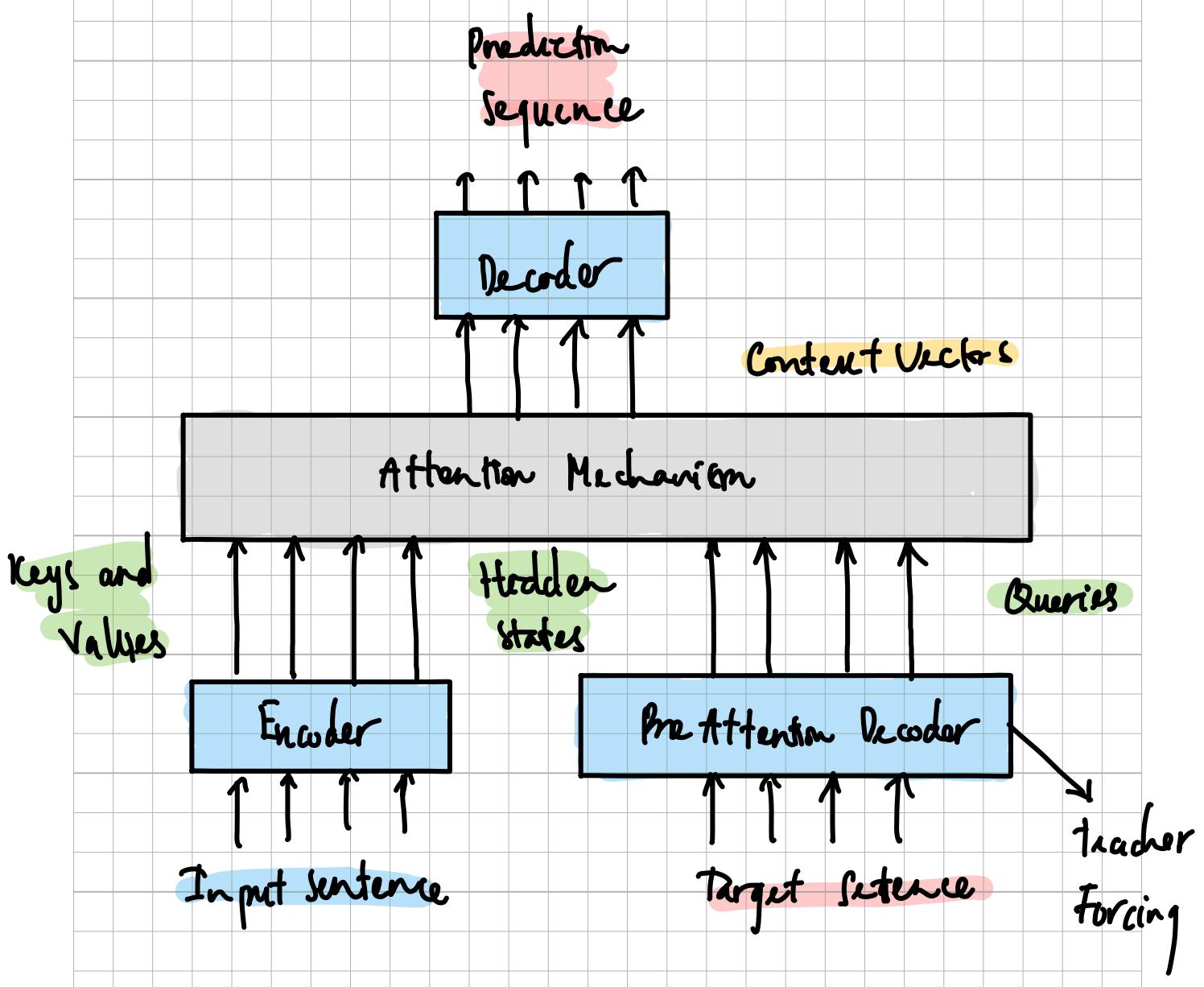
(Naïve model)

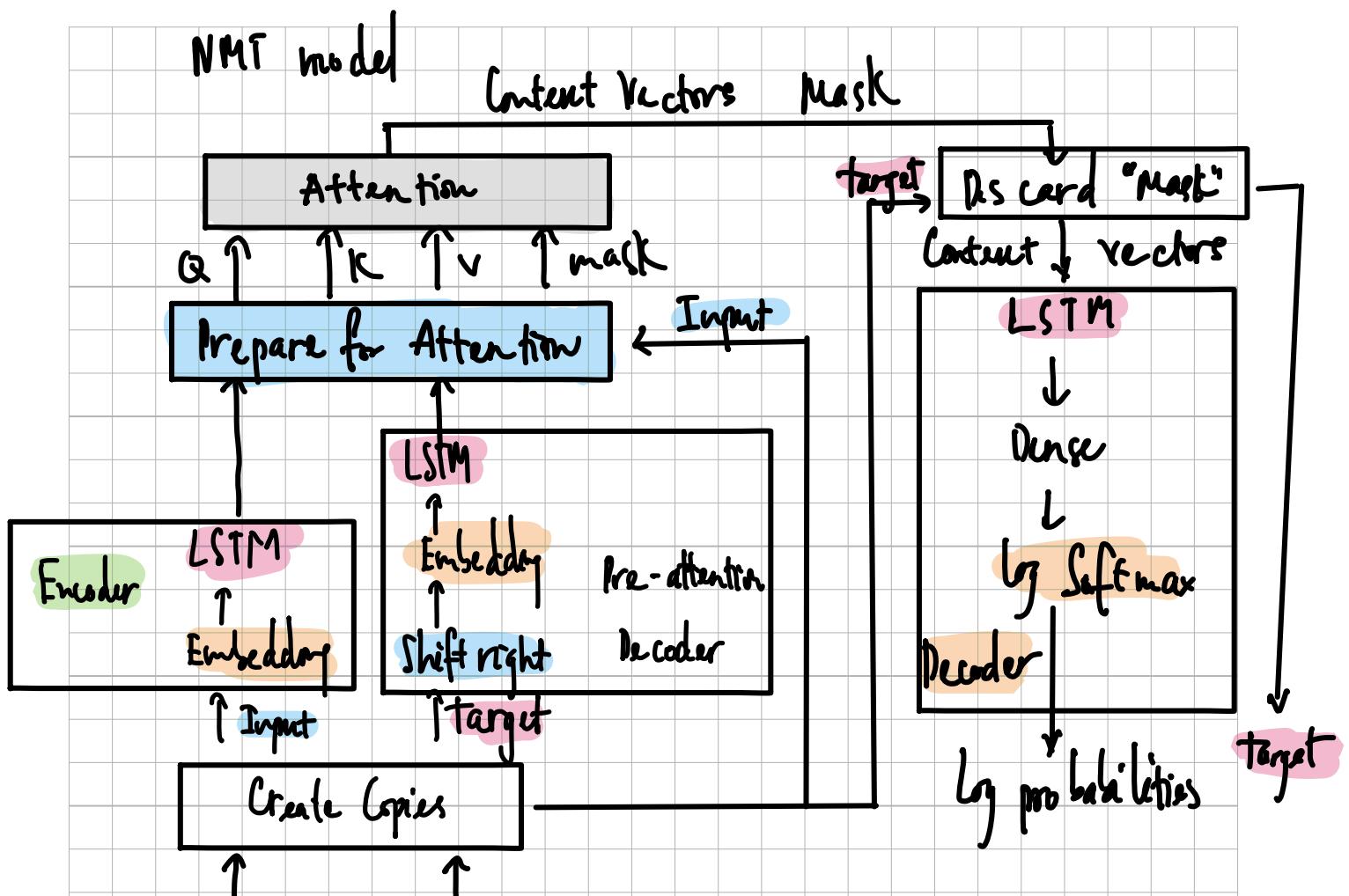
Irl → Correct sequence of words as input (shifted right)

Improves training performance

slowly using decoder outputs

b) NMT with Attention





+ BLEU score (Modified)

bitingual Evaluation Under way

Compare candidate translation to reference

the better (human) translations

↔ 1

Does not consider Semantic meaning,
sentence structure

Precision

+) ROUGE-N score : recall-oriented

$$F_1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$\rightarrow F_1 = \frac{\text{BLEU} \times \text{ROUGE}-N}{\text{BLEU} + \text{ROUGE}-N}$$

only account N-grams

Random Greedy B:1 matching

+ Sampling and Decoding

Temperature control more or less randomness in predictions

lower temperature, more confident, conservative which
higher temperature, more excited, random output

+ Beam search

best sentences over fixed - window size

beam width

→ Calculate probability of multiple possible sequences
at each step

beam width B: # sequences you keep

Until all B most probable sequences end with <EOS>

beam search decoding

Penalizes long sequences, so you should normalize by the sentence length

Computationally expensive and
consumes a lot of memory

+ Minimum Bayes Risk (MBR)

$$\underset{\substack{E \\ \downarrow}}{\operatorname{argmax}} \frac{1}{n} \sum_{\substack{E' \\ \text{compare}}} \text{ROUGE}(E, E')$$

rouge score

find the candidate with every between pair of candidates
that maximizes other candidate

(13) Text summarization

+) Transformers vs. RNNs

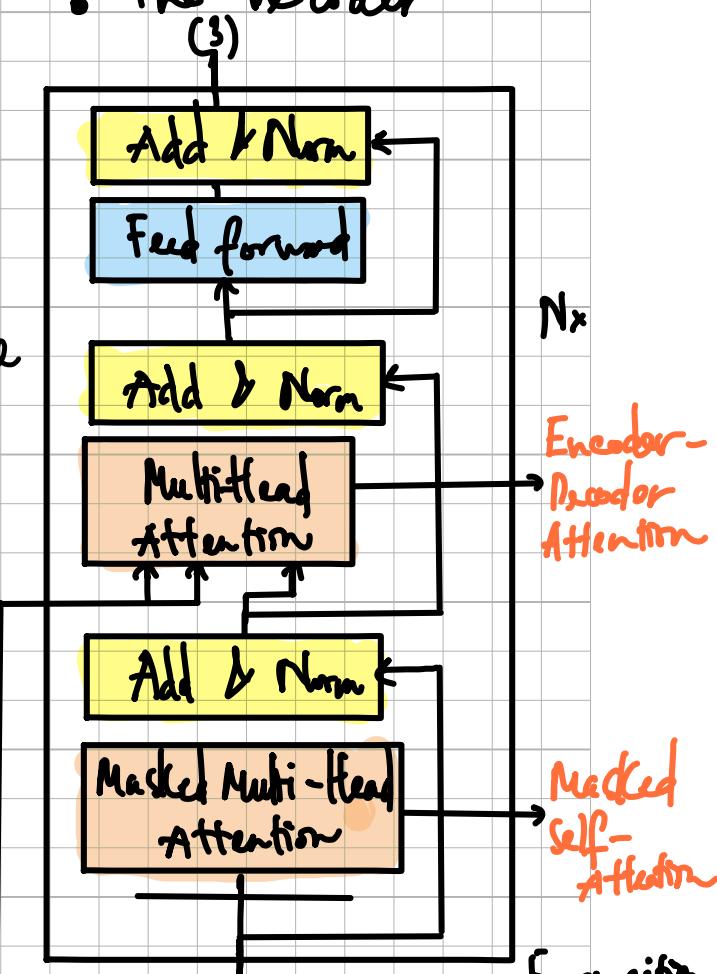
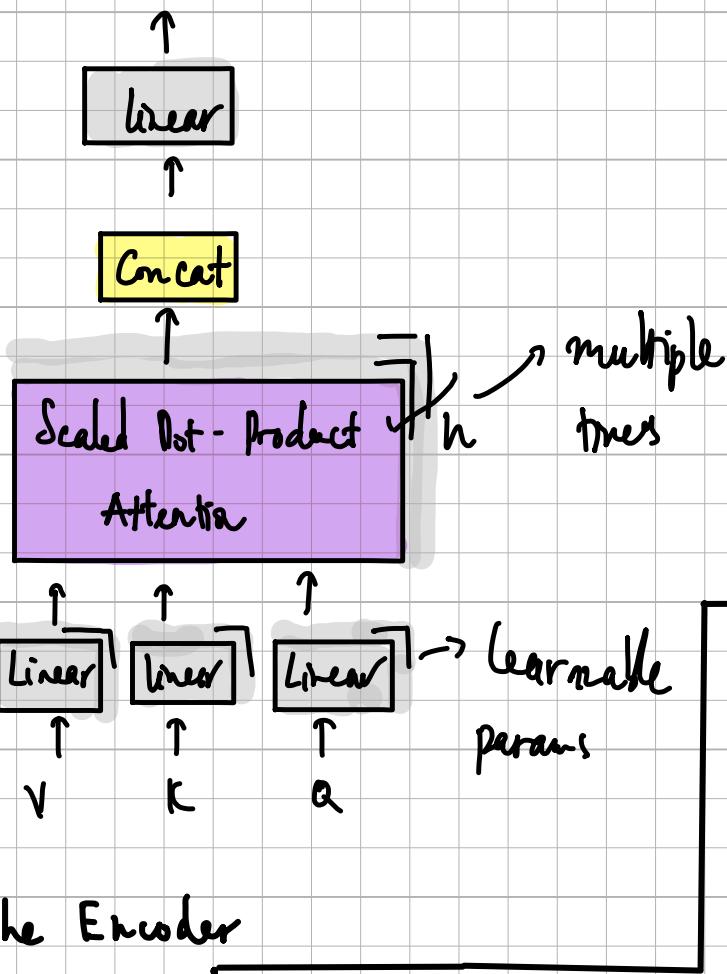
(no parallel)

2d7 Google

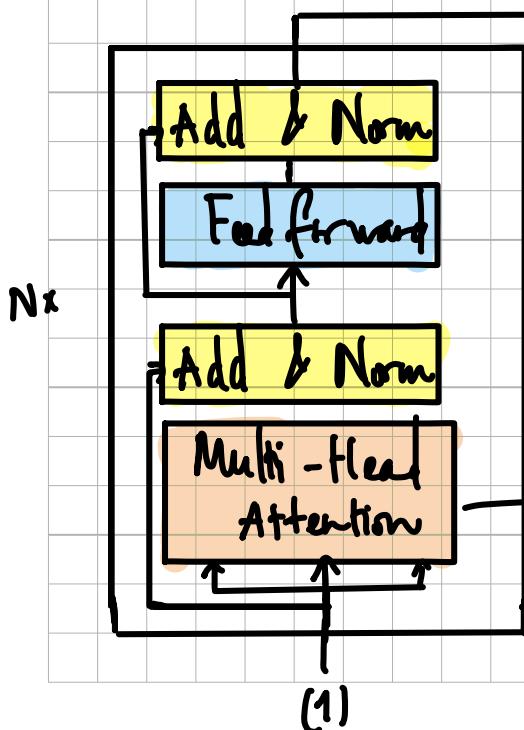
↳ sequential computation,

loss of information, vanishing gradients

- Multi - Head Attention



- The Encoder

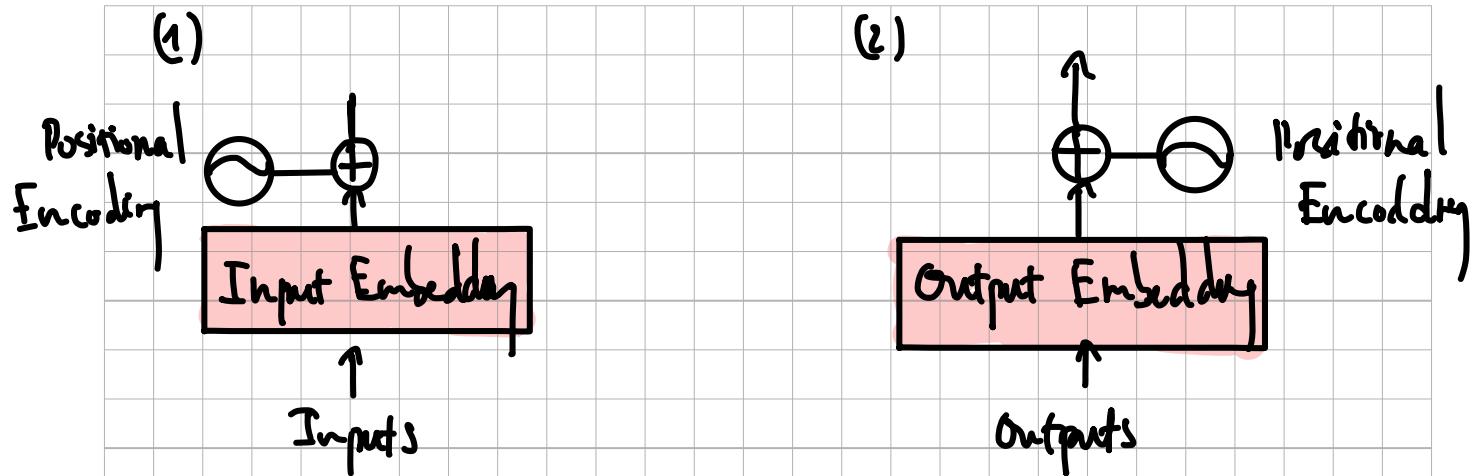


Provides contextual representation of each item in the input sequence

• Positional Encoding

Self Attention

→ Every item in the input attends to every other item in the sequence



(3)

Probabilities

Softmax

Linear

+ Transformer Applications

Text summarization Translation

Auto complete

Chatbots

NER

Other NLP tasks

Q & A

OpenAI

2018

GPT2

Google AI language 2018

Google

2019

BERT

T₅: Text-to-Text Transfer

Transformer

(multitask)

Translation

Classification (CoLA sentence)

QA

Regression (SST-2)

Summarization

Bi-directional

Encoder Representations

from Transformers

Generative

Pretraining

Transformer

+1 Scaled and Dot - Product Attention (Heart and Soul)

$$\text{softmax} \left(\frac{\mathbf{Q} \mathbf{C}^T}{\sqrt{d_K}} \right) \mathbf{V} = \text{context vectors for each query} (\# \text{queries} \times \text{size of value vector})$$

w[2,3] weight assigned to third key, second query

\mathbf{Q} : stack of rows embedding $m \times \text{size of embedding}$ (generally the same)

\mathbf{K}, \mathbf{V} : same number of rows \times size of embedding

Attention Math

+1 Masked Self Attention

• Encoder- Decoder
Attention

Queries from one sentence,
Keys and values from another

• Self - Attention

$\mathbf{Q}, \mathbf{K}, \mathbf{V}$ comes from the same sentence
Meaning of each word within the sentence

• Masked Self Attention

$\mathbf{Q}, \mathbf{K}, \mathbf{V}$ come from the same sentence

Queries don't attend to future positions

$$\text{Softmax} \left(\frac{\mathbf{Q} \mathbf{K}^T}{\sqrt{d_K}} + \begin{pmatrix} 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{pmatrix} \right) \rightarrow -\infty$$

+1) Multi-headed Attention

sets = # heads

$d_{\text{model}} = \text{embedding size}$

$$\underbrace{W_i^Q \quad W_i^K}_{(d_{\text{model}} \times \frac{\text{heads}}{d_{\text{model}}})} \quad \underbrace{W_i^V}_{\downarrow d_{\text{model}}} \quad W^O$$

usual choice of dimensions

$$d_K = d_V = d_{\text{model}} / h$$

Attention

$$d_V$$

→ Concat $h \times d_V$

→ Context vectors for each query
 d_{model}

Causal Attention \approx Masked Attention

$$\text{PE}_{(\text{pos}, 2i)} = \sin\left(\frac{\text{pos} \cdot i}{10000^{2i/d}}\right)$$

$$\text{PE}_{(\text{pos}, 2i+1)} = \cos\left(\frac{\text{pos} \cdot i}{10000^{2i/d}}\right)$$

+2) Transformer Decoder

feed forward (ReLU)

+3) Transformer for Summarization

Technical details for data processing

loss weight: 0 until the first <EOS> and then
is on the start of summary

(0 → 0, 2 0, 5)

$$J = -\frac{1}{n} \sum_j^m \sum_i^L y_j^i \log \hat{y}_j^i$$

i: batch elements

j: over summary

Model Input

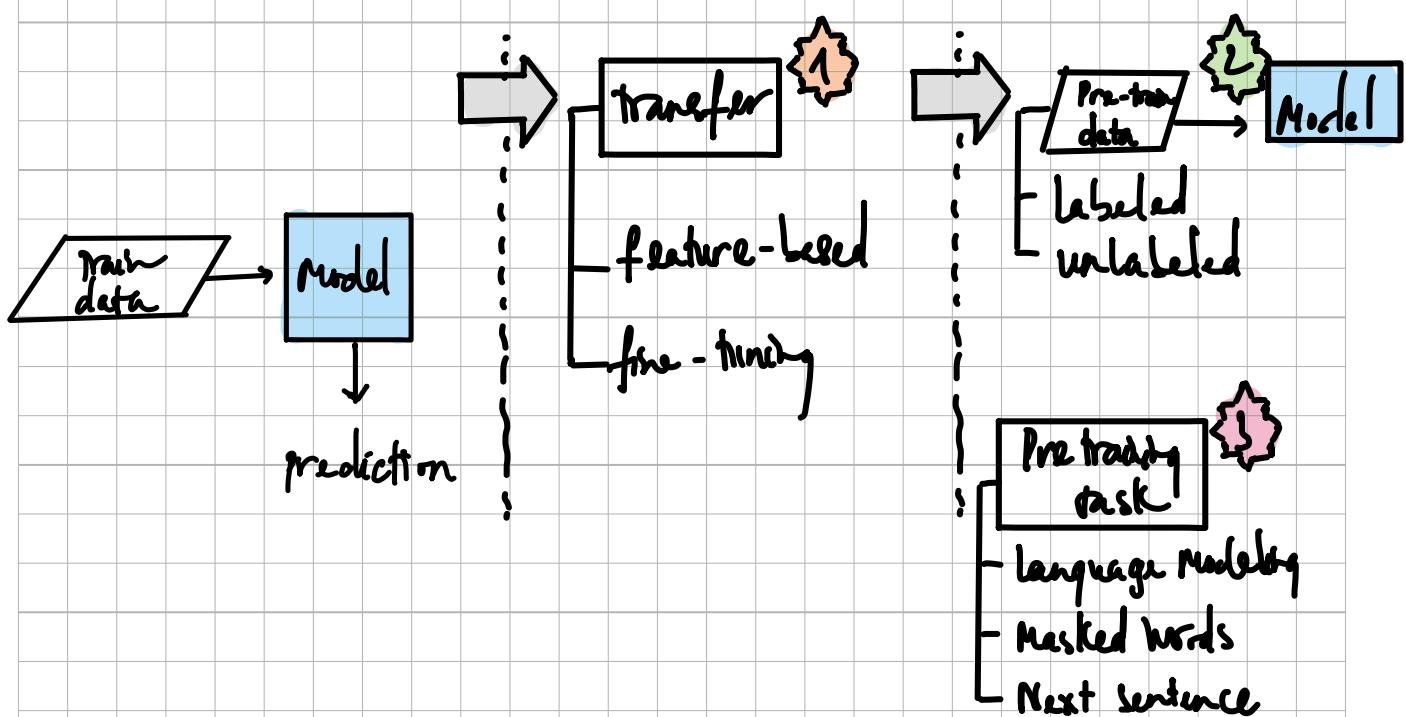
[Article] <EOS> [summary] <EOS>
(, tokenized input

(14) Question Answering

Content-based

Closed - box

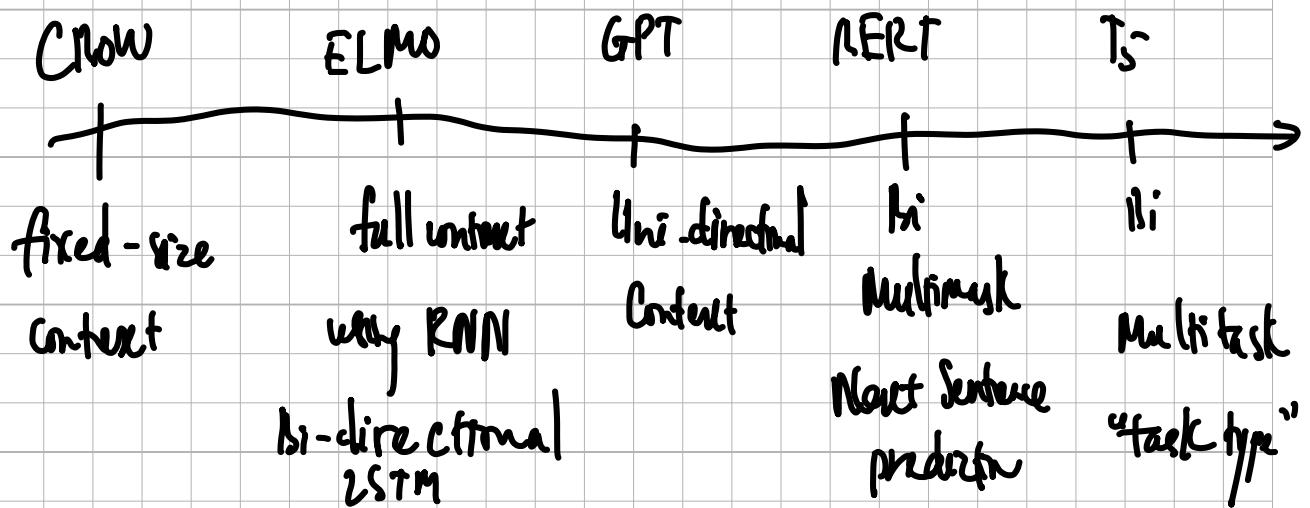
+ Transfer learning with full transformer



i) General purpose learning (1)

Self-supervised tasks (1)
(U)

ii) ELMO, GPT, BERT, TS



+) Bidirectional Encoder Representations from Transformers

BERT

• Makes use of transfer learning / pretraining

• A multi layer bi-transformer

• Positional embeddings

• BERT - Base

12 layers (12 transformer blocks)

12 attention heads

110 million params

processes 15% words

During pretraining, the model is trained on unlabeled data over different pre-training tasks

- Choose 15% of tokens at random

- mark 80% of the time, replace with random

- Keep as is 10%, token 10%

- multiple masked spans in a sentence

- Next Sentence prediction

for finetuning, initialized Bert with pretrained params, using labeled data from the downstream tasks.

+) BERT objective

combine word embedding, sentence embeddings, positional embedding

Formalizing the Input

Input

[CLS] ... [SEP] ... [SEP]

- + Token Embeddings
- + Segment Embeddings
- + Position Embeddings

C
↓
NSP

T_1, \dots, T_n [sep] T'_1, \dots, T'_m
mask LM mask LM

o Objective 1 :

Multi mask LM

Loss: Cross Entropy Loss

o Objective 2 :

Next Sentence Prediction

Loss: Binary Loss

(Do they follow or Not)

t) Fine tuning BERT

Sentence A
text

Sentence B
∅

Question

Passage

Hypothesis

Premise

Sentence

Editor

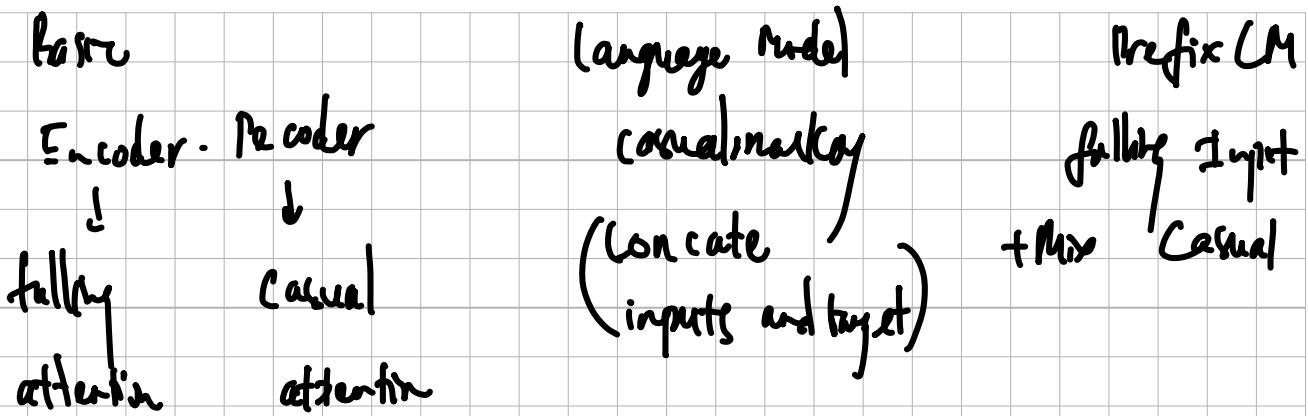
Sentence

Paraphrase

Article

Summary

t) T5 : Transfer learning + Masked language modeling
< >



+ Multi-task training strategy

- Predict entailment, contradiction, or neutral

- Winograd scheme + *

Training data strategies:

Examples - proportional mix, temperature scaled mix

Equal mix

Fine-tuning method:

- All params

Adapter layers

Gradual unfreezing

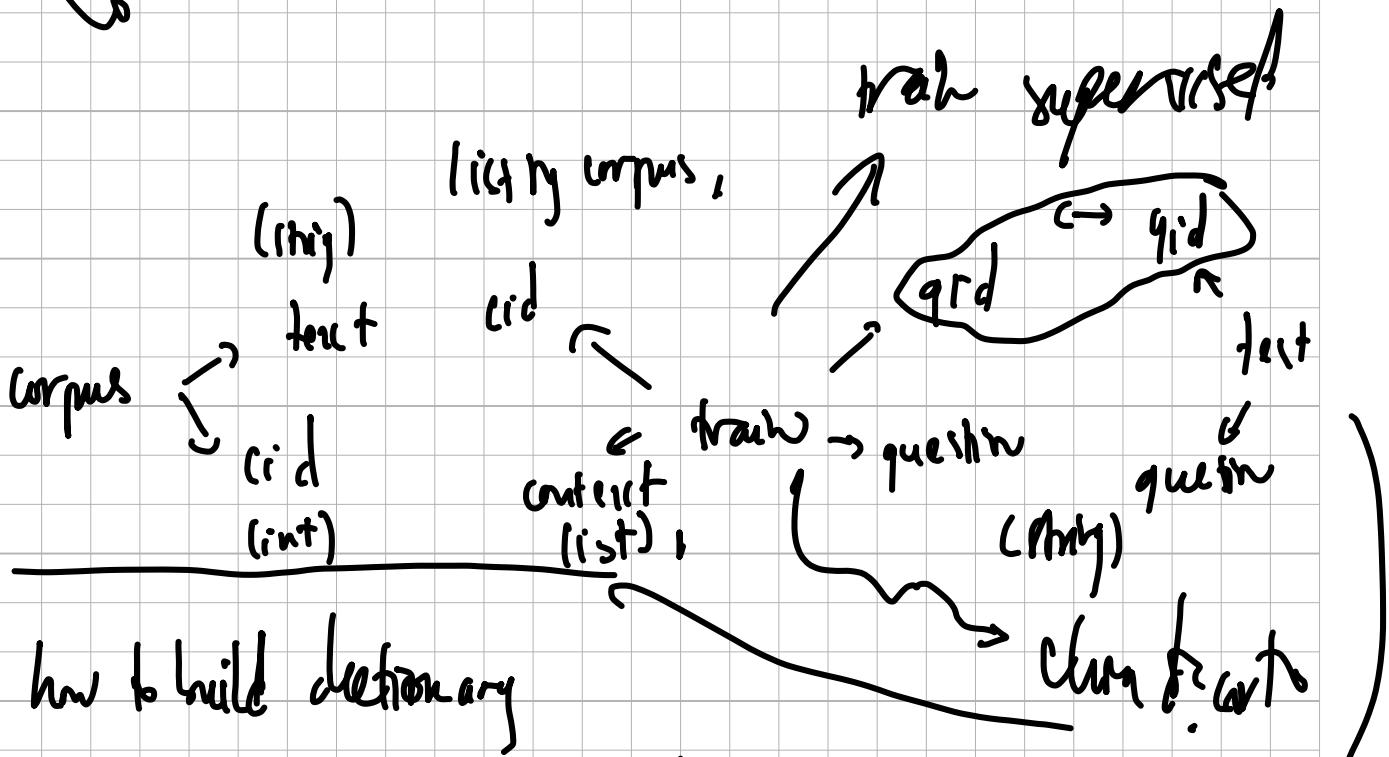
+ GLUE benchmark

+ triggering Facc : using transformers

Pipeline

1. Pre-processing your inputs
2. Running the model
3. Post-processing the outputs

↳



Search qid → uid top k

Can be o α

gid?

Sort freq qid

Similarity?

Summarize list and context
and question
→ top k cod

photent py-vn corelp fair seg

Vietnamese - bi-encodes