

Natural Language Processing

* Probabilistic language models

→ conditional probability

→ The Chain Rule

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

→ Markov Assumption: reduced context window

→ start symbol <S> and stop symbol </S>

* Unigram, Bigram, Trigram

$$\hookrightarrow P(w_1) \times P(w_2) \times \dots \times P(w_n) : P(w_i) = \text{count}(w_i) / \text{Total words}$$

$$Bi: P(w_2 | w_1) \times \dots \times P(w_n | w_{n-1}) : P(w_2 | w_1) = \text{count}(w_1, w_2) / \text{count}(w_1)$$

$$Tri: P(w_3 | w_1, w_2) \times \dots \times P(w_n | w_{n-2}, w_{n-1}) : P(w_3 | w_1, w_2) = \text{count}(w_1, w_2, w_3) / \text{count}(w_1, w_2)$$

* we count stop symbol as a word in our dictionary.

Example

$$P(I \text{ am}) = \frac{\text{count}(I \text{ am})}{\text{count}(I)} = \frac{2}{3}$$

$$P(\text{am Sam}) = \frac{\text{count}(\text{am Sam})}{\text{count}(\text{am})} = \frac{1}{2}$$

* Issues with probabilistic models:

→ zero probability: unseen data leading to zero overall

→ underflow: small prob leading to near zero prob

* underflow: take logs of prob multiplications (solution)

* zero prob: smoothing is the solution

* The Shannon Visualization method

→ next word prediction method

* The intuition of smoothing is to maintain the true distribution by giving weight which is borrowed from seen data to unseen data.

* Laplace smoothing: pretend we saw each word one more time than we did.

Train set

{ <S> I ate an apple today </S>

Example

<S> I at a banana </S>

<S> strawberry is sweet </S>

test set

{ <S> I ate an orange today </S>

$$\text{Laplace} = \text{count(word)} + 1 \quad N = 15 : \text{total words}$$

$$N + V \quad V = 11 : \text{unique vocab}$$

$$\text{so for word "I": } \frac{\text{count}(I) + 1}{N + V} = \frac{2 + 1}{15 + 11}$$

Backoff, Interpolation.

* Linear Interpolation:

$$P(w_n | w_{n-1}, w_{n-2}) = \lambda_1 P(w_n | w_{n-1}, w_{n-2}) + \lambda_2 P(w_n | w_{n-1}) + \lambda_3 P(w_n)$$
$$\therefore \sum \lambda_i = 1$$

* Set lambdas using held-out corpus

→ Training data, held out data, best data

* Unknown word '`<unk>`'

* Stupid backoff: doesn't make a distribution

`<s> I saw van </s>`

`<s> saw a cat </s>`

`<s> he ate an orange </s>`

`<s> I saw an orange </s>`

Bigrams with linear interpolation

$$P(<s> I) = \lambda_1 P(<s> I) + \lambda_2 P(I)$$

$$P(I saw) = \lambda_1 P(I saw) + \lambda_2 P(saw)$$

$$P(saw an) = \lambda_1 P(saw an) + \lambda_2 P(an)$$

$$P(an orange) = \lambda_1 P(an orange) + \lambda_2 (orange)$$

↑ ↑
bigrams unigrams

* Intrinsic, Extrinsic Evaluation

$$\begin{aligned} * \text{Perplexity} &= P(w_1, w_2, \dots, w_n)^{1/N} \\ &= \sqrt[N]{\frac{1}{P(w_1, w_2, \dots, w_n)}} \\ &= 2^{-1/N} \leq \lg P(w_i) \end{aligned}$$

* Text preprocessing

* **Text Classification**

→ aspect based sentiment analysis

* **Classification Methods:**

→ Hand-coded rules

→ supervised machine learning

* **Naive Bayes Classification**

→ Negation dealing.

→ Naive Bayes works like unigram, words are independent

→ Bag of words representation ← prior prob of class

$$P(c|d) = \frac{P(d|c) P(c)}{P(d)} \text{ where 'd' is the document and 'c' is the class}$$

$$P(\text{pos}|d) = P(d|\text{pos}) P(\text{pos}) / P(d)$$

$$P(\text{neg}|d) = P(d|\text{neg}) P(\text{neg}) / P(d)$$

* $P(d)$ from denominator can be removed in cases where the documents are same.

* **Multinomial naive bayes independence assumptions**

$$P(x_1, x_2, \dots, x_n | c) = P(x_1 | c) \cdot P(x_2 | c) \cdots P(x_n | c)$$

$$\text{maximum likelihood estimate} = P(\text{class}) = \frac{\text{document}(C=c_i)}{N_{\text{doc}}}$$

$$P(\text{word}/\text{class}) = \frac{\text{count of word in class}}{\text{count of all words in class}}$$

$$P(c) = \frac{3}{4} \Rightarrow$$

$$P(\text{chinese}/j) = \frac{1}{3}$$

$$P(\text{chinese}/c) = \frac{5}{8} =$$

$$P(\text{Beijing}/c) = \frac{1}{8}$$

$$P(\text{Shanghai}) = \frac{1}{8}$$

$$P(\text{Macau}) = \frac{1}{8}$$

without smoothing

$$P(d|c) = \frac{3}{4} \left(\frac{5}{8} \times \frac{5}{8} \times \frac{5}{8} \times \frac{0}{8} \times \frac{0}{8} \right) = 0$$

$$P(d|j) = \frac{1}{4} \left(\frac{1}{3} \times \frac{1}{3} \times \frac{1}{3} \times \frac{1}{3} \times \frac{1}{3} \right) = 0.001$$

with smoothing (Laplace).

$$P(d|c) = \frac{3}{4} \left(\frac{(5+1)^3}{8+6} \times \frac{0+1}{8+6} \times \frac{0+1}{8+6} \right) = 0.000301$$

$$P(d|j) = \frac{1}{4} \left(\frac{(1+1)^3}{3+6} \times \frac{(1+1)}{3+6} \times \frac{(1+1)}{3+6} \right) = 0.0000148.$$

* Binarized Naive Bayes / Boolean

→ prioritizes the presence or absence of words more than the count of the words.

→ the data is changed such that there are no duplicates in each document.

Text Classification Evaluation

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{precision} = \frac{TP}{TP + FP}$$

$$\text{recall} = \frac{TP}{TP + FN}$$

$$\text{true recall} = \frac{TP}{TP + FN}$$

how many were actually true

out of all that were actually true how many did the model predict correctly?

		P	N
		TP	TN
		FP	FN
ground truth		gold pos	gold neg
predicted	system pos	true pos	false pos
	system neg	false neg	true neg

Confusion matrix

$$A = \frac{25+20}{25+10+40+20}$$

Example

		TP	FP		Accuracy = 0.45
C1:		25	45	Prec = $\frac{25}{25+45} = 0.357$	
		10	20	Rec = $\frac{25}{25+10} = 0.714$	
		FN	TN		

C2:	10	30	P = $\frac{10}{10+30} = 0.25$
	25	35	R = $\frac{10}{10+25} = 0.2857$
			Accuracy = $\frac{10+35}{10+30+25+35} = 0.45$

* F-Score: $\frac{2PR}{P+R}$ $\therefore P = \text{precision}, R = \text{recall}$

↳ harmonic mean of precision and recall

Q: Why do we use harmonic mean instead of arithmetic mean?

A: In harmonic mean, if one value is extremely small, it has a huge impact on the result which helps in identifying a heavily biased model.

Macro and Micro Average

Example: $\text{OK} : P = \frac{95}{95+10} \Rightarrow \frac{95}{105} \text{ for a } 6 \times 6 \text{ confusion matrix}$

$$R = \frac{95}{95+1+13+1} = \frac{95}{110} \text{ micro average}$$

Macroaveraging:

Compute the precision and recall of all classes and then take the average, all classes ^{don't} have same weightage because the number or size of class will affect the weightage of class in total averaging.

Example

	# of obj	Precision
C1:	1000	0.1
C2:	30	0.8
C3:	500	0.01
C4:	100	0.04

$$\text{MacroAvg} = P = \frac{0.1 + 0.8 + 0.01 + 0.04}{4}$$

$$P = 0.2375.$$

* in micro, a pooled confusion matrix is formed

pooled:

✓ for example on slides.

$$\begin{array}{c|c} 268 & 99 \\ \hline 99 & 635 \end{array} = \frac{268}{268+99} = 0.73$$

* micro-averaging assigns equal weights to each "object" not each class.

K-Fold validation

* divide the validation dataset in 'k' and introduced each in every iteration for validation.

* vectors of words, dimensions are vocabulary

Embeddings

→ Term-document matrix (word x doc)

→ Term-context matrix (word x word)

	apple	orange	to	and	of	cat	dog	and	the	window=2
pe	0	2	2	3	2	1	2	0	0	
dog	0	0	0	2	2	1	0	2	1	

* euclidean, cosine

* Euclidean Distance

$$\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

* Cosine Similarity

$$\cos(\vec{v}, \vec{w}) = \frac{\vec{v} \cdot \vec{w}}{\|\vec{v}\| \|\vec{w}\|} = \frac{\sum_{i=1}^N v_i w_i}{\sqrt{\sum_{i=1}^N v_i^2} \sqrt{\sum_{i=1}^N w_i^2}}$$

$$\cos(\text{april, digital}) = 0 + 0 + 0 / = 0$$

$$\cos(\text{april, info}) = 1 + 0 + 0 / \sqrt{1^2} \sqrt{1^2 + 2^2} = 0.33$$

$$\cos(\text{dig, info}) = 0 + 6 + 2 / \sqrt{1^2 + 2^2} \sqrt{1^2 + 6^2 + 1^2} = 0.58$$

- * Raw frequency is a bad representation as it does not take into consideration the difference in the length of the documents.

so, we normalize our vectors

→ length normalization → divide word count by doc length

→ \log_{10} normalization → $1 + \log_{10}$ (word count)

* IDF : inverse document frequency

IDF = total number of documents

(number of document frequency (w.r.t word))

→ log transformed $\log_{10} \left(\frac{N}{df_w} \right)$

$$TF-IDF = TF \times IDF$$

Example:

	good	cat	food	doc-frequency of words
d1	60	0	100	good
d2	100	50	0	cat
d3	0	20	50	food
d4	0	1	1	

$$TF-IDF = (60) \log_{10} \left(\frac{4}{60} \right) + 1 \times \log_{10} \left(\frac{1}{1} \right)$$

$$TF-IDF_{d1} = 1 + \log_{10}(60) \times \left[1 + \log_{10} \frac{10^6}{100,000} \right]$$

$$= 4.556 + 0 + \left[1 + \log_{10}(100) \times \frac{10^6}{100,000} \right]$$

$$= 4.556 + 1 = 18.556$$

* TF-IDF vector with context windows.

Neural Language Model

→ Sampling

→ top K sampling.

* window based neural language model

Cross entropy:

$$y \log y + (1-y) \log (1-y)$$

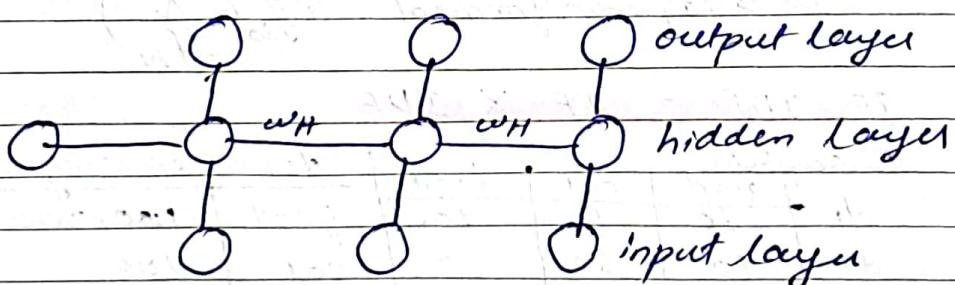
drawbacks of neural language model:

→ slow training on more parameters.

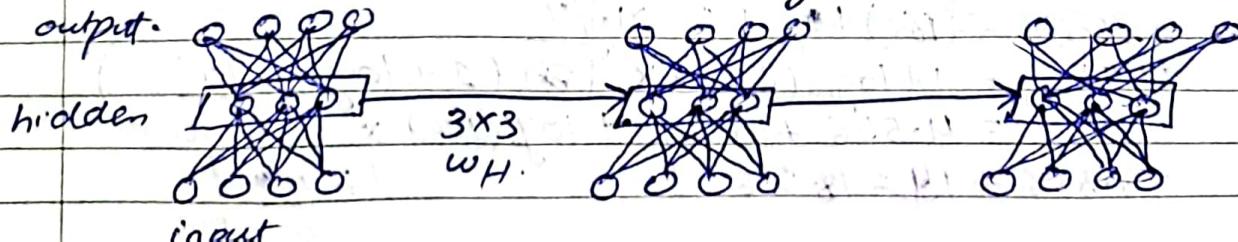
→ ~~no ability to~~

→ no symmetry in how the inputs are processed i.e. the nearby neighbours & sequence of words

Recurrent Neural Network



* In RNNs, the hidden layers are connected output.



$$y_3 = w_y h_3$$

$$h_3 = w_x(x_3) + w_h(h_2)$$

$$h_2 = w_x(x_2) + w_h(h_1)$$

$$h_1 = w_x(x_1) + w_h(h_0)$$

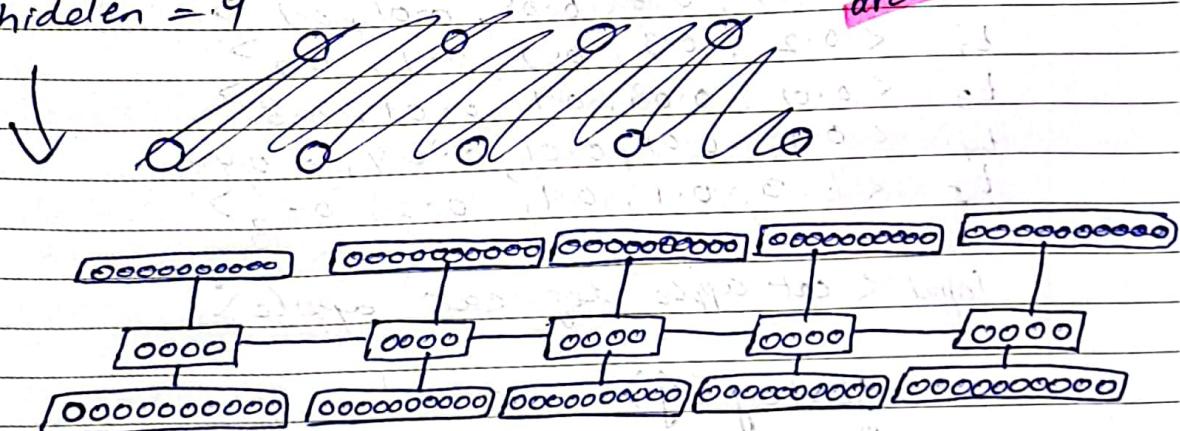
→ this shows how the first word ~~to~~ x_1 would contribute in the last output y_3 .

→ vanishing gradient problem

vocab = 10

input word sequence = 5

hidden = 4



advantages over neural network:

→ gives more weight to closely context words ~~to~~ and less to far away neighbours

→ number of input words doesn't affect the number of weights thus doesn't slow down training

disadvantages:

→ vanishing gradient problem

cross-entropy

$$= \underbrace{y \log y^{\dagger}}_{y=1} + \underbrace{(1-y) \log (1-y)}_{y=0}$$

* cross-entropy measures how similar or dissimilar the ground truth & the predicted distributions.

<dog, cat, won, apple, eat>

$$t_1 <0.04, 0.8, 0.05, 0.01, 0.1>$$

$$t_2 <0.2, 0.3, 0.2, 0.1, 0.2>$$

$$t_3 <0.01, 0.08, 0.1, 0.01, 0.8>$$

$$t_4 <0.9, 0.01, 0.01, 0.04, 0.04>$$

$$t_5 <0.3, 0.1, 0.1, 0.2, 0.3>$$

Input <cat apple dog eat won>

$$\frac{-1}{n} \sum \log(y^{\dagger})$$

$$\frac{-1}{n} [\log(0.8) + \log(0.1) + \log(0.01) + \log(0.04) \\ + \log(0.2)]$$

gradient descent:

$$\text{chain rule: } \frac{d(L)}{d(u)} = \frac{d(J)}{d(v)} \times \frac{d(v)}{d(u)}$$

$$J(w) = -\log y^{\dagger}$$

$$\frac{\partial J(w)}{\partial w_2}$$

$$\frac{d(L)}{d(x_1)} = \frac{d(J)}{d(v)} \times \frac{d(v)}{d(u)} \times \frac{d(u)}{d(x_1)}$$

$y_t = \text{softmax}(z_t)$
 $z_t = w_y h_t$
 $h_t = w_{hx} \cdot x_t + w_{hh} \cdot h_{t-1}$
 $\frac{\partial L_T}{\partial w_{hx}} = \frac{\partial L_T}{\partial y_T} \cdot \frac{\partial y_T}{\partial z_T} \cdot \frac{\partial z_T}{\partial h_T} \cdot \frac{\partial h_T}{\partial w_{hx}}$

but h_{t-1} also contributes to the

- * vanishing gradient problem
in long input sequences the gradient becomes so small that it doesn't contribute significant change in the far away weights.

* LSTM.

has gates.

the forget gates solves the vanishing gradient problem as the gradient doesn't diminish on all nodes, only on contributing nodes.

Machine Translation.

- * parallel corpora

* Neural Machine Translation (NMT)

encoder \rightarrow decoder.

\hookrightarrow the error & loss & weight updates only apply to the decoder part.

Beam search

Beam Search: select some outputs of a decoder layer & pass it forward to generate some output sequences then compute the sentence probability of these generated sequences, select the best.

Breadth First: generate all possible sentences & compute their sentence probabilities choose the best.

$$\text{prob. of } \text{S} = \text{prob. of } \text{S} \cdot \text{prob. of } \text{S} = 0.16$$

$$\text{prob. of } \text{S} \cdot \text{prob. of } \text{S} = 0.064$$

out of 256 possible states, go to

the next condition, first two paddings of sequence and according length of sequence, the next frame of search will be based on previous step's student's answer.

Step 1

Step 2

Answers are based upon beam search, since it is not able to consider all possibilities, so there is a window function, which filters out those no useful and invalid sequences.

Final answer is based upon beam search.

Final answer is based upon beam search.

Final answer is based upon beam search.