

23/01/24NLP

Task Contamination → doesn't work on data that haven't been seen before in LLMs.

Edit Distance:candidate words: giraffe graf grail

{ Spell-Checking }

Coreference: Harvard President Claudine
→ Howard University " "

They are coreference because they are directed towards a particular word.

The apple was red

2\$ 2\$

8\$ 4\$

1. The apple was big and red

2. The fish was blue.

cost: 4 \$

cost: 8 \$

By calculating the cost we can find which sentence is more similar.

Operation and cost

Insertion \rightarrow 2\$

Deletion \rightarrow 2\$

Substitution \rightarrow

What is the minimum edit distance to make.

INTENTION \rightarrow EXECUTION

$\downarrow \downarrow \downarrow$
 * EXECUTION

d s s i s

cost: 8

L^E_AVANSHTEIN

insertion — 1\$

deletion — 1\$

substitution — ~~d + i~~
 $= 2\$$

* Edit Distances

1. Levenshten
2. Hamming
3. Jaro-winkler

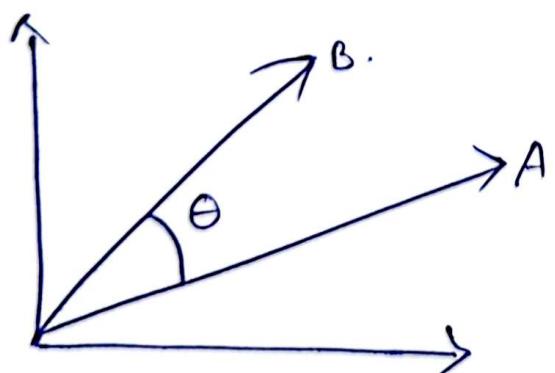
* Token Based:

1. Cosine - similarity
2. Jaccard - Index
3. Tversky - Index

* Sequence Based:

1. LCS.
2. Ratcliff - Obershelp Similarity

Cosine-similarity



The more similar two sentences are, the more it will coincide making θ smaller and smaller.

cosine similarity (doc1, doc2)

$$= \frac{\vec{A} \cdot \vec{B}}{|A| \cdot |B|}$$

Doc 1: The cat in the hat.

Doc 2: A black cat in a black hat.

Unique tokens:

	The	Cat	in	Hat	a	black
Document 1	2	1	1	1	0	0
Doc-2	0	1	1	1	2	2

$$\text{vector } A = \begin{bmatrix} 2 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

$$B = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 2 \\ 2 \end{bmatrix}$$

$$\text{similarity} = \frac{\vec{A} \cdot \vec{B}}{|A| |B|} = \frac{(2 \times 0) + (1 \times 1) + (1 \times 1) + (1 \times 1) + (0 \times 2) + (0 \times 2)}{\sqrt{7} \cdot \sqrt{11}}$$

$$= 0.342$$

Perpendicular vector (no similarity) close to 0

Disadvantages

→ can't detect semantics of sentence.

Doctor ~~at~~ came after patient died
Doctor can before patient died

(Opposing sentence).

→ No negative values.

→ Sensitive towards length of the vector. More length will have more sparsity.

Sparsity: A matrix/vector having lots of zeros

Jaccard Similarity:

$$JS(doc_1, doc_2) = \frac{doc_1 \cap doc_2}{doc_1 \cup doc_2}$$

→ I wake up early

→ I go to Canada for Phd.

} no

} similarity

Cat, mat

= $\{ \text{'the, cat; in, the, mat}' \} \cap \{ \text{'A, black, cat, in, a, black, mat;} \}$.

$\{ \text{'The, cat, in, mat, a, black}' \}$.

= $\frac{3}{6} \cong 50\%$.

Lemmatization: Root of word

throw threw thrown

N-gram Language Models.

Please turn in your homework

Converting the sentence to bi-grams we get:

Bi-grams = Please turn
turn in.
in your
your homework

tri-grams = Please turn in
turn in your
in, your homework

Skip bi-grams = Please in
turn your
in homework

Language Model: Assigns probabilistic values to language data to do a particular task. Like GPT does to generate/predict the next word.

Markov Model: To predict something you don't have to look too far. So:-



$$\rightarrow P(W_n | W_{n-1}) = \frac{Q(W_{n-1}, W_n)}{Q(W_{n-1}, W)}$$

We can predict the next word by just looking at the next word.

$\rightarrow P(I | <s>)$ [what is the probability the 'I' is the starting of the string].

$$= \frac{Q(<s> I)}{C(<s>)} = 2/3$$

$<s>$ I am Palmer </s>
 $<s>$ Palmer I am </s>
 $<s>$ I do not like eggs </s>

$$\rightarrow P(am | I) = 2/3$$

$$\rightarrow P(</s> | Abear) = \frac{C(Abear </s>)}{C(Abear)}$$

$$= 1/2 \dots ,$$

Longer grams contain more content.

Shorter grams often fails to capture the context.

Rule-Based: Hardcoded probability assignment.

Feature-Based - Captures all probability automatically

$\langle s \rangle$ I want Chinese Food $\langle /s \rangle$

$$P(\langle i | \langle s \rangle) = 0.25$$

$$P(\langle /s \rangle | \text{food}) = 0.68$$

We are given probabilities of bi-gram data:-

$$P(\langle s \rangle \text{ I want Chinese Food } \langle /s \rangle)$$

$$= P(i | \langle s \rangle) P(\text{want} | i) P(\text{Chinese} | \text{want})$$

$$P(\text{Food} | \text{chinese}) P(\langle /s \rangle | \text{Food})$$

$$= 0.25 \times 0.33 \times 0.065 \times 0.52 \times 0.68$$

$$= 0.0002$$

$$\Rightarrow \log(P_1 \times P_2 \times P_3) = \log P_1 + \log P_2 + \log P_3$$

$$\Rightarrow P_1 + P_2 + P_3 = e(\text{value.})$$

30/01/24

NLP

* Extrinsic Evaluation :

e.g. Prediction models like gmail reply suggsts.

→ How well a model performs on real-world tasks.

* Intrinsic Evaluation :

e.g. Used in a simulated environment which means

Test-Based scenario to evaluate, for assessment.

e.g. Understanding Syntactic Structure .

Q: How can we say if a Language Model fits test-data better?

→ The probability of the prior language model is better.

Q: What happens if I ~~feed~~ evaluate my language model with 'Seen' data

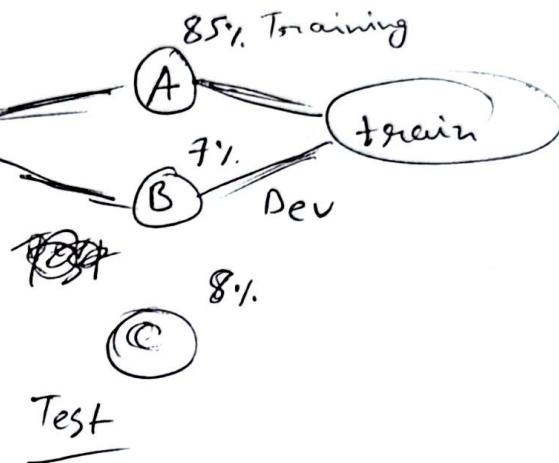
→ Overfitting: Memorizes the data instead of generalizing. Biasness increases. Performs bad in Training set. But performs well in

Underfitting: Low Biasness. Performs bad in the Test, Training and validation set.

- Why shouldn't we train with 100% data?

⇒ Doesn't represent real-world scenario.

Train-test-split



EVALUATION METRIC - PERPLEXITY

→ Used to evaluate how well a model performs generative tasks.

$$\text{Perplexity, } \text{PPL}(w) = N \sqrt{\frac{1}{P(w_1, w_2, \dots, w_N)}}$$

$$= 2^{-\frac{1}{N} \sum_{i=1}^N \log_2 P(x_i)}.$$

"The cat is on the mat" (Test Sample).

$$P("The") = 0.1 \quad P("Cat") = 0.2 \quad ("is") = 0.15 \\ ("on") = 0.1 \quad ("the") = 0.2 \quad ("mat") = 0.25.$$

$$PPL = 2^{-\frac{1}{6}(\log_2(0.1) + \dots + \log_2(0.25))} \\ = 2^{-\frac{1}{6} \times (-15.02)} = 2^{2.5} \approx 5.66.$$

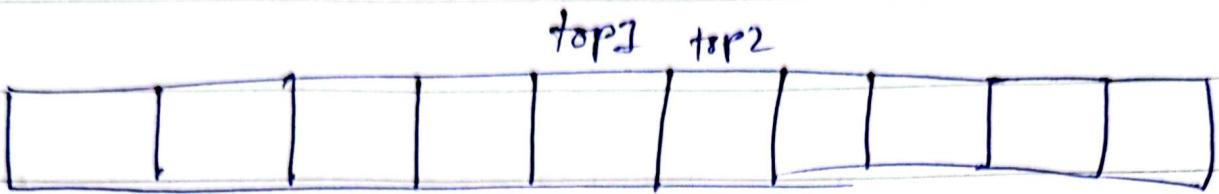
\Rightarrow There are six branches in the training corpus where the above probabilities will found.

The lower the Perplexity, the better the language model in Predicting.

(

Algo Lab

t = 0



06-02-24

NLP

Disadvantage of a large N-gram Model.

Smoothing / Discounting: Take some from the probabilities of other values and add to counts that are zero.

The simplest Smoothing Method is Laplace method.

We increase all the probabilities by 1.

So, all 0s become 1.

$$P = \frac{\text{count}(w, v) + 1}{\text{all count of bigrams} + 1} \quad | \quad \begin{aligned} \text{count}(i|i) &= 5 \\ P(i|i) &= 5/2533 \end{aligned}$$

$$\text{Lap. } P(i|i) = \frac{6}{2534} \quad \begin{matrix} a \\ b \\ c \end{matrix} \quad \text{my(2)}$$

I want food lunch.

$$P = P(I | \langle ss \rangle) \times P(\text{want} | I) \times P(\text{food} | \text{want}) \times P(\text{lunch} | \text{food}) \\ P(\langle ss \rangle | \text{lunch}).$$

A variant of Laplace smoothing is K-smoothing.

Add-K. # How do you get K? $\frac{\text{count}(wv) + K}{\text{all counts of bigram} + K}$.

Backoff, Interpolation

If a probability is not given at let say, Bigram we backoff and move to the Unigram.

$P(\text{lunch} | \text{food}) = 0$, & in bigram. So, we backoff to Unigram and take that probability.

Interpolation:

$$P(w_n | w_{n-1}, w_{n-2}) = \lambda_1 P(w_n) + \lambda_2 (\text{Bigram } p) + \lambda_3 (\text{Trigram } p).$$

$$\therefore \sum^n = \lambda p.$$

Oahr () Oahr

Knesser-Ney

I can't see without my reading

francisco

Pcontinuation $\propto \{v : CC(vw) > 0\}$

techniques vs before ws

$$\frac{P_c(\text{glasses})}{P_c(\text{Francisco})} = \frac{2}{1} \quad \left| \begin{array}{l} \text{Higher Probability of} \\ \text{occurring in different} \\ \text{contents} \end{array} \right.$$

Sentiment Classification

* Binary Cls: Email spam or not.

* Multinomial Cls: Discrete classification.

* Ordinal

* Multiclass C

* MultiLevel Classif

Bag of Words: Doesn't store the order of words.

$$\hat{c} = \underset{\text{argmax}}{\text{argmax}} \ P(c|d) \quad \text{document}$$

$$\underset{\text{argmax}}{\frac{c \in C}{\text{argmax}}} = \frac{P(c) \times P(d|c)}{P(d)}$$

$$= \underset{c \in C}{\text{argmax}} \frac{P(c)}{I} \times \frac{P(d|c)}{\text{(Likelihood)}}$$

If from a document,

Prior: Probability of documents taken from class. ($P(c)$)

$P(d|c)$: Probability that the document is already in class.

$$P(d|C) = P(f_1 \cdot f_2 \cdot f_3 \cdots f_n | C)$$

→ all possible combination of features.

This is very expensive = 100^c_{100} .

To simplify it, we take a Naive approach:-

$$P(f_1|C) \cdot P(f_2|C) \cdot P(f_3|C) \cdots P(f_n|C).$$

So,

$$\hat{C} = \underset{C \in \mathcal{C}}{\operatorname{argmax}} \prod_{f \in F} P(f|C) = \underset{C \in \mathcal{C}}{\operatorname{argmax}} P(C) \prod_{i \in \text{Position}} P(w_i|C)$$
$$= \underset{C \in \mathcal{C}}{\operatorname{argmax}} \log P(C) + \sum_{i \in \text{Position}} \log(P(w_i|C))$$

They are called as linear classifiers - aka Generative Classifier.

Spam: 2, Normal: 3, Total: 5. $n_{\text{doc}} = 2/5$ $n_{\text{norm}} = 3/5$.

Likelihood: $P(w_i|C) = \frac{\text{count}(w_i, c)}{\text{vocabulary} \times \sum_{(w \in V)} \text{count}(w, c) + 1}$ Laplace Smoothing.

$$= \frac{\text{count}(w_i, c) + 1}{\sum_{(w \in V)} \text{count}(w, c) + |\mathcal{V}|}$$

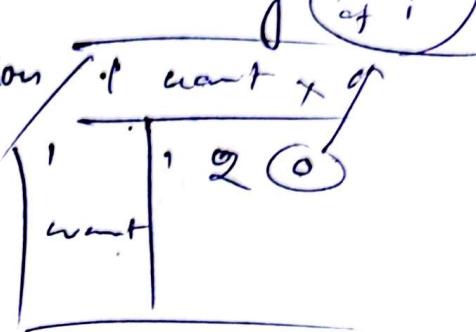
Ignore the words that aren't in bag of words.

When to do Laplace:

Ignore word that never occurred in training (occurred of i)

Corpus.

$$S = \text{Win} \oplus \text{free vacation}$$



$$P(\text{spam}) = \frac{3}{6} \quad |v| =$$

$$P(\text{ham}) = \frac{3}{6}$$

$$P(\text{win}|\text{spam}) = \frac{1+1}{10+16}$$

$$P(\text{free}|\text{spam}) = \frac{2+1}{10+16}$$

$$P(\text{vacation}|\text{spam}) = \frac{0+1}{10+16}$$

$$P(\text{win}|\text{ham}) = \frac{0+1}{8+16}$$

$$P(\text{free}|\text{ham}) = \frac{0+1}{8+16}$$

$$P(\text{vacation}|\text{ham}) = \frac{1+1}{8+16}$$

$$P(c) \cdot \prod P(w_i|c)$$

$$P(\text{spam}) \times P(S|\text{spam}) = \frac{3}{6} \times \frac{2+3+1}{26^3} \cong 0.000171$$

$$P(\text{Ham}) \times P(S|\text{Ham}) = \frac{3}{6} \times \frac{1+1+2}{24^3} \approx 0.000072.$$

I really like this book
I don't like this book.

Bag of words cannot handle negation expression

Modification:

Prepending: Add NOT before each token until punctuation when encountered or a negative word (don't, isn't).

I don't NOT-like NOT-this NOT-book

LWC - Linguistic Inquiry

MPQA

+ve : admirable, beautiful,
(0.9)

-ve: Catastrophe, envious
(0.4)

Evaluation Metric

	Span	Han
Sends Han to Span %.	99%.	1%.
time but 99% time it		
classifier Span -		
so, Accuracy 1 99%.		

But this Accuracy although number-wise good, it's unuseful and not beneficial.

		gold (+ve)	gold (-ve)	spam
		System (+ve)	true (+ve)	precision = $\frac{tp}{tp + fp}$
		System (-ve)	false (-ve)	true (-ve)
			recall = $\frac{tp}{tp + fn}$	accuracy = $\frac{tp + tn}{tp + tn + fp + fn}$

100

Spam (-ve) \rightarrow 99

Han. (+ve) \rightarrow 1

acc = 99%

rec = 0

false (+ve) = 0 [no spam classified as Han].

true Positive = 0 [Han not correctly identified].

$$\underline{f\text{-measure}} = \frac{(\tilde{\beta}+1)PR}{\beta^2 P + R} \rightarrow \frac{\text{precision}}{\text{recall}}$$

$(\beta > 1) \rightarrow R$
 $(\beta < 1) \rightarrow P$

NLP

Feb 13

Confusion Matrix

	U	S	n	
U	89	10	1	U = urgent
S	5	60	50	S = spam
n	3	36	252	n = non-spam

$$P_U = \frac{89}{89+10+1}, \quad S_U = \frac{8}{8+5+3}$$

But how do we measure

* macro-averaging

* micro-averaging



Class 1 : urgent

8	11
8	340

 P_{n_1}

~~anything~~ Every cell (row & column) wise are added other than 1st cell because they are false negatives.

$$P_m = \frac{8}{8+11}$$

268	99
99	635

micro

Class 2

60	55
90	212

 P_{n_2}

Class 3

200	33
51	83

 P_{n_3}

$$\text{macro-avg precision} = \frac{P_{n_1} + P_{n_2} + P_{n_3}}{3} = 0.60.$$

$$\text{micro avg precision} = \frac{268}{268+99} = 0.71$$

Report , when all classes are equally important , use macro .

When we want to report a particular class , use micro .

Given a class is dog, what are the features it has:

$$P(\text{cl}d) = \hat{C} = P(C) \cdot \underbrace{P(d|c)}_{\text{likelihood}} \quad [\text{General Classifier}]$$

Logistic Regression Classifier:

Discriminates between two class. Dogs wearing belt vs no cats wearing belt. So, puts a lot of weight on belt and discriminates. finds the Discriminative Feature

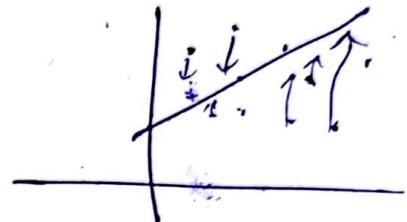
Logistic Regr. Functions:

- i) Feature Representation
- ii) Classification Function
- iii) Objective Function.
- iv) Optimization Function.

Linear Regression :

Mean Square Error =

$$\frac{\sum_{i=1}^n (y_i - y_i^P)^2}{n}$$



Mean Absolute Error =

$$\frac{\sum |y_i - y_i^P|}{n}$$



loss

Hill Climbing Method (Local Search). $y = mx + b$

for random m & b to minimize MSE.

$(\overset{m, b}{0, 2}) \rightarrow (1, 2)$ [Less MSE?] .
[go forward]

Which Loss Function should we use?

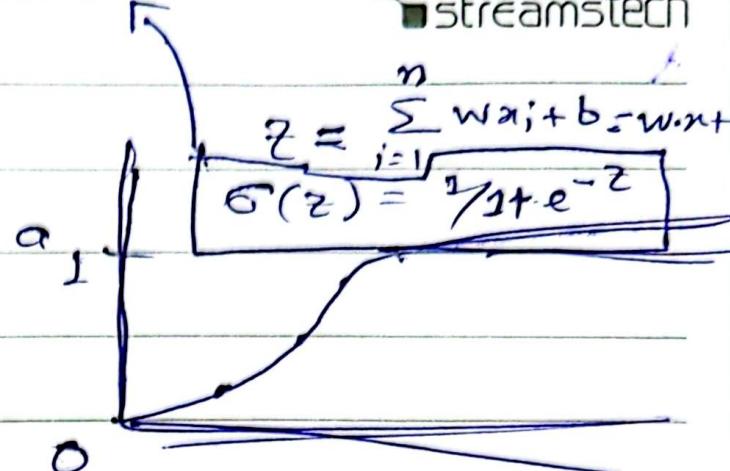
Use equal weights to all sample $\rightarrow L_2$

Ignore Outliers (Sakit) $\rightarrow L_1$.

SIGMOID

$$y = \frac{a}{1 + e^{-(wx+b)}}$$

Sigmoid Function.



Do not rise above 'a'.

0

Decide w if the point is closer to 1 or 0
(Binary Classification).

Now let's change the terminology:-

w = weight, b = bias

$$\therefore y = \frac{w}{1 + e^{(-wx+b)}}$$

Calculation:

$$x = [3, 2, 1, 3, 0, 4.19].$$

$$w = [2.5, -5, -1.2, 0.5, 2, 0.7] \quad [b = 0.1].$$

$$P(+ | \text{doc}) = \sigma(w \cdot x + b).$$

$$= \sigma(0.833) = \frac{1}{1 + e^{-0.833}} \approx 0.7$$

$$P(- | \text{doc}) = 1 - \sigma(z) = 1 - 0.7 = 0.3$$

NLP
15/02/24

Previously we saw a doc with only one 6 features and drew a regression line. But when the features are very large in number we create a matrix instead of manually calculating the probabilities.

$$\hat{y} = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(m)} \end{bmatrix} \begin{bmatrix} w \end{bmatrix} + \begin{bmatrix} b \end{bmatrix}$$

\hat{y} [m x 1] x [f x 1] w [f x 1] b [1 x 1]

$$\hat{y} = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_f^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \dots & x_f^{(2)} \\ \vdots & \vdots & & \vdots \\ x_1^{(m)} & x_2^{(m)} & \dots & x_f^{(m)} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_f \end{bmatrix} + \begin{bmatrix} b^{(1)} \\ b^{(2)} \\ \vdots \\ b^{(m)} \end{bmatrix}$$

$$\boxed{\hat{y} = w \cdot x + b}$$

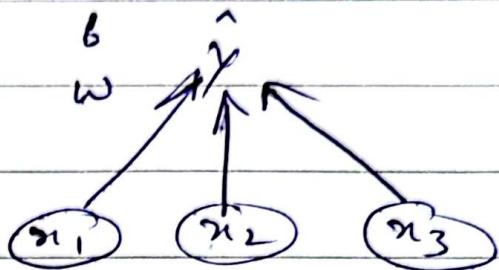
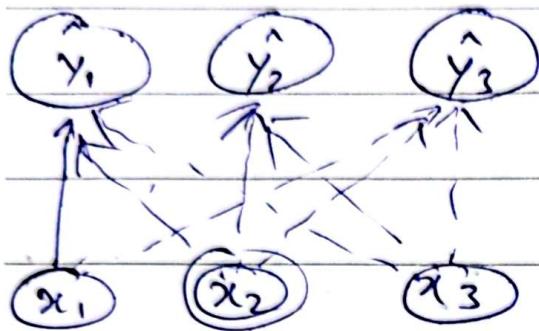
m = number of document
 f = features

* Sigmoid Function *
 (Binary Classifier)

- Classifies only one class.
- Scknets. Software

for multiple classes,

for single class



(Activation Functions),

Sigmoid: $z = w \cdot x + b$, $z = \frac{e^z}{1+e^{-z}}$

Softmax: $\hat{y}_i = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}}, z = [0.7, 0.15, 0.02 \dots]$

$L(\hat{y}, y)$ = How much \hat{y} differs from the true y

Bernoulli Distro (when Probability is either 0 and 1).

$$P(y|x) = \hat{y}^y (1-\hat{y})^{1-y} \quad [\text{Gold Level}] = 0.$$

$$\log p(y|x) = y \log \hat{y} + (1-y) \cdot \log(1-\hat{y}) .$$

$$L_{CE}(\hat{y}, y) = - [\hat{y} \log \hat{y} + (1-\hat{y}) \log(1-\hat{y})] .$$

(cross entropy)

$$= - [y \log \sigma(wx+b) + (1-y) \log(1-\sigma(wx+b))]$$

When we can differentiate it, we can optimize it.

-infinity to 0 is the range. A good classifier will have value close to 0. A bad classifier will go towards -infinity.

Gradient Descent: Slope of a straight line or curve.

Loss Function: Curve \rightarrow when value is 0.

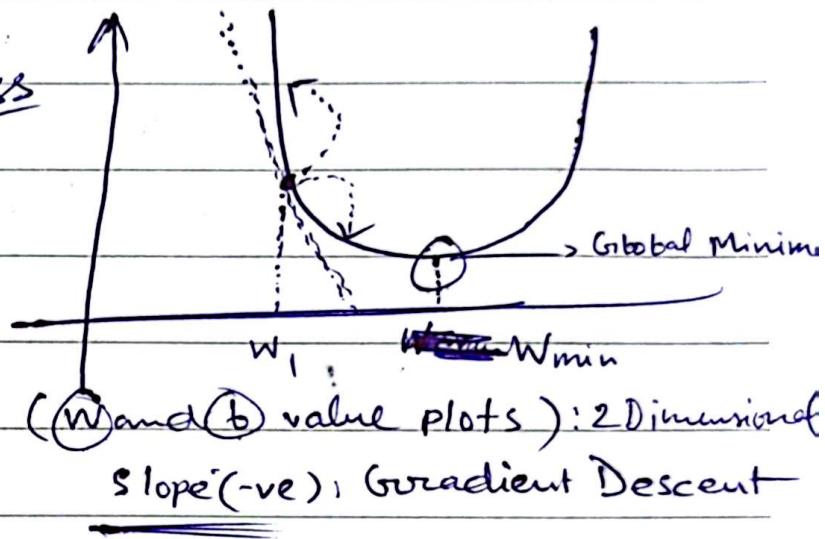
$$\hat{\theta} = \underset{\theta}{\operatorname{arg\,min}} \frac{1}{m} \sum_{i=1}^m L_{CE}(f(x^{(i)}, \theta), y^{(i)})$$

(w, b) ↴
 (w, b)

To minimize go downward

Global
Minima?

STUCK!



(eta)

$$w^{t+1} = w^t - \eta \frac{d}{dw} (L(f(x, w), y))$$

η = learning rate

means how much downward

you'll go to reach global
minima.

$$x = [x_1, x_2], \quad w_1 = w_2 = b = 0, \quad \eta = 0.1.$$

$$x_1 = (+ve) \text{ lexicous} = 3.$$

$$x_2 = (-ve) \text{ lexicous} = 2.$$

$$\theta^{t+1} = \theta^t - \eta \nabla L(f(u_i, \theta), y_i) \quad \left| \begin{array}{l} \sigma(\theta) = \frac{1}{1 + e^{-\theta}} \end{array} \right.$$

$$\nabla_{w,b} L = \begin{bmatrix} \frac{\partial}{\partial w_1} L_{CE}(\hat{y}, y) \\ \frac{\partial}{\partial w_2} L_{CE}(\hat{y}, y) \\ \frac{\partial}{\partial b} L_{CE}(\hat{y}, y) \end{bmatrix} = \begin{bmatrix} [\sigma(w \cdot x + b) - y] x_1 \\ [\sigma(w \cdot x + b) - y] x_2 \\ [\sigma(w \cdot x + b) - y] y \end{bmatrix} = \begin{bmatrix} -0.5x_1 \\ -0.5x_2 \\ -0.5 \end{bmatrix} = \begin{bmatrix} -1.5 \\ -1.0 \\ -0.5 \end{bmatrix}$$

$$L_{CE} = -[\gamma \log(\hat{y}) + (1-\gamma) \log(1-\hat{y})].$$

$$\frac{\partial}{\partial w_j} L_{CE} = (\hat{y} - y) x_j \rightarrow \sigma(w \cdot x + b)$$

$$\therefore \theta + \epsilon I = \begin{bmatrix} w_1 = 0 \\ w_2 = 0 \\ b = 6 \end{bmatrix} - 0.1 \begin{bmatrix} -1.5 \\ -1.0 \\ -0.5 \end{bmatrix} = \begin{bmatrix} 0.15 \\ 0.1 \\ 0.05 \end{bmatrix}$$

Stochastic Gradient Descent .



20/02/24 - NLP.

$$\hat{\Theta} = \arg \min_{\Theta} \frac{1}{m} \sum_{i=1}^m L_{CE}(f(x^{(i)}; \Theta), y^{(i)}) + R(\Theta).$$

L_2 Norm:

L_1 & L_2 Norm. \rightarrow to solve Overfitting

L_2 Norm: Uses Euclidian distance

for L_2 , $R(\Theta) = \alpha \sum_{j=1}^n \Theta_j^2$

L_1 Norm: $R(\Theta) = \alpha \sum_{i=1}^n |\Theta_i|$

- Makes the bigger weights 0
- Makes a Sparse Vector (Breaks larger weights).
-

Embedding \rightarrow Express Words as vectors in vector space based on their meaning.

- Previously we used Frequencies to classify features
- Representation Learning \rightarrow Time Consuming / Bad.

Solution:

Vector Semantics: Projecting words into a multi-dimensional space based on their meaning.

Mouse — Rodent

|

Pointer Device

→

Synonyms:

when two Lemmas have the same meaning.

(Mouse) → Mice

↓

→ root word is known as

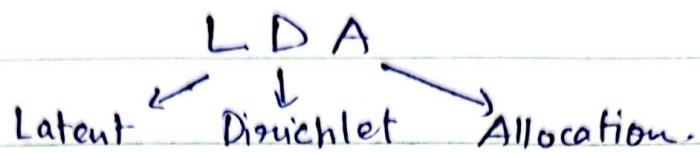
Lemmas

→ word forms

word
senses

Principle of Contrast: Similar words are not always substitutable.

Related Words;



In a certain interval, which words occur more frequently? When we have a document and don't know the topic, we can use LDA.

* Connotations : \rightarrow i) Valence ii) Arousal iii) Dominance

Excited (-, -, -)

Projecting a word in multi-dimensional space

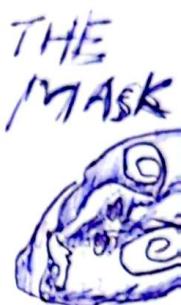
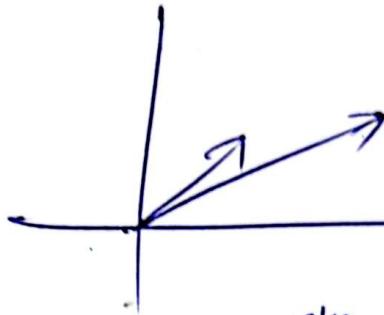
22 Feb 22

NLP

- Similar words contain similar documents.
- Similar documents contain similar words.

Cons:

Depending too much on frequency
STILL.



$$\text{idf} = N/\text{df} \quad \boxed{\text{tf-idf} = \text{tf} * \text{idf}}$$

$$= \log_{10}(N/\text{df})$$

vector space.

If a term appears in all of the documents then
 $\text{idf} = 0$

Positive

Pointwise Mutual Information

NLP

27 Feb.

Embedding: Dense Vectors (not too sparse - data having value zero),

TF-IDF and PPMI are too much frequency dependant which result in a Sparse Vector Output.

We use 'Word2vec' model to actually learn the weights and represent the data as vectors.

Focus Word; Content words;

C BOW

→ We know the focus word, context words but we want to predict a focus word.

Skip - Ngram:

→ We know the focus words and want to predict Content words

→ Chances of error/fault is more

Purpose

→ Skip N-gram can be used to classify the rare words more.

Negative Sampling

To find Content word we take tuples of word.

[Pad]	Data	is	processed	into	information	[Pad]
	c	f	c			
	c	f	c			

→ The centre word is the focus word.

→ Content Window: 1 in here, before: 1 word,
after: 1 word.

Here every word appears as a content word and focus word at least once.

Negative Sampling: Soln to disaster caused by softmax Function.

$$\log p(c|w; \theta) = \frac{\exp^{v_c \cdot v_w}}{\sum_{c' \neq c} \exp^{v_{c'} \cdot v_w}}$$

Scale Unigram with $P = \frac{c^\alpha}{\sum c^\alpha}$ $\alpha = 0.75$.

So, frequent words are scaled down $(3)^{0.75}$

non-frequent words are scaled up $(0.001)^{3/4} > 0.001$

What If the words do occur in the context window, then is a 1 or 0.

Inherent Variability:

- Create Multiple embeddings
- Bootstrap & Average to find value.