# Language Modeling .

Word embeddings
in a nutshell



WORDS GET A VECTOR, DOCS GET
A VECTOR

EVERYTHING GETS A VECTOR

word
vectors

context

Word
embeddings

When you penalize your Natural
Language Generation model for
large sentence lengths

Me think, why waste time say lot
word, when few word do trick.

do a little text search and nobody
bats an eye

do a little
NLP and everybody loses their minds

quickmeme.com

# Recap:

* Corpus $\rightarrow$ Collection of documents.

* Documents $\rightarrow$ Collection of one or more sentences.

* Sentences $\rightarrow$ Collection of words.

* Converting text datasets to vectors:

    (i) Document $\rightarrow$ Vector 1,2.

    (ii) Word $\rightarrow$ Vector. 3,4.

* Methods for vectorization:

    (i) Bag-of-words

    (ii) TF-IDF

    (iii) Continuous Bag-of-words

    (iv) Skip-gram.

## Agenda:

* **Language modeling** — teaching a computer how to form sentences.

* **Approaches**: What could be done here?

* **Techniques**: Unigram, Bigram, Tri-gram, n-gram.

* **Core Concept**: Conditional Probability.

* But first... business case!!

# Naïve approaches - How to solve the problem?

* Conditional prob. → Naïve Bayes.

* Word 2 Vec → Cosine similarity.
   ↳ problem: similar words, maybe not so useful.

* Take CBoW → modify to produce the last word, we could have a NN. solution.

## Conditional Prob.

A, B

$P(A)$, $P(B)$.

$P(A \cap B)$

$\checkmark P(A|B) = \dfrac{\overset{\downarrow}{P}(A \cap B)}{P(B)}$

$P(B|A) = \dfrac{P(A|B) \cdot P(B)}{P(A)}$

A and B are independent

if $P(A|B) = P(A)$

How to apply conditional probability to predict the next word?

$$\begin{cases} w_1 \\ w_2 \\ w_3 \\ \vdots \\ \vdots \\ w_n. \end{cases}$$

* The cat is ___ .

$w_1 \quad w_2 \quad w_3 \quad w_4.$

$w_4$ is nothing but the word which maximizes the following probability

$$w_4 = \arg\max_i P(w_i \mid w_1 = \{the, \; w_2 = cat, \; w_3 = is \; )$$

# Problems with probabilities:

1] $P(w_R \mid w_{k-1} \, w_{k-2} \, w_{k-3} \cdots w_1)$

$$= \frac{P(w_R \cap w_{k-1} \cap w_{k-2} \cap \cdots)}{P(w_{k-1} \cap w_{k-2} \cap \cdots)}$$

low values. $\longrightarrow$ log probabilities.

2] $P(w) = \underline{0} \longrightarrow$ happens when we have not seen $w$ before.

$\hookrightarrow$ Laplace smoothing.

## Joint probability

$$P(w_1 \cap w_2 \cap w_3 \cap \dots \cap w_n) \leftarrow \frac{O}{\uparrow}?$$

↳ "What is the probability of seeing

→ $\underline{w_1, w_2, w_3, \dots, w_n}$ in exactly this
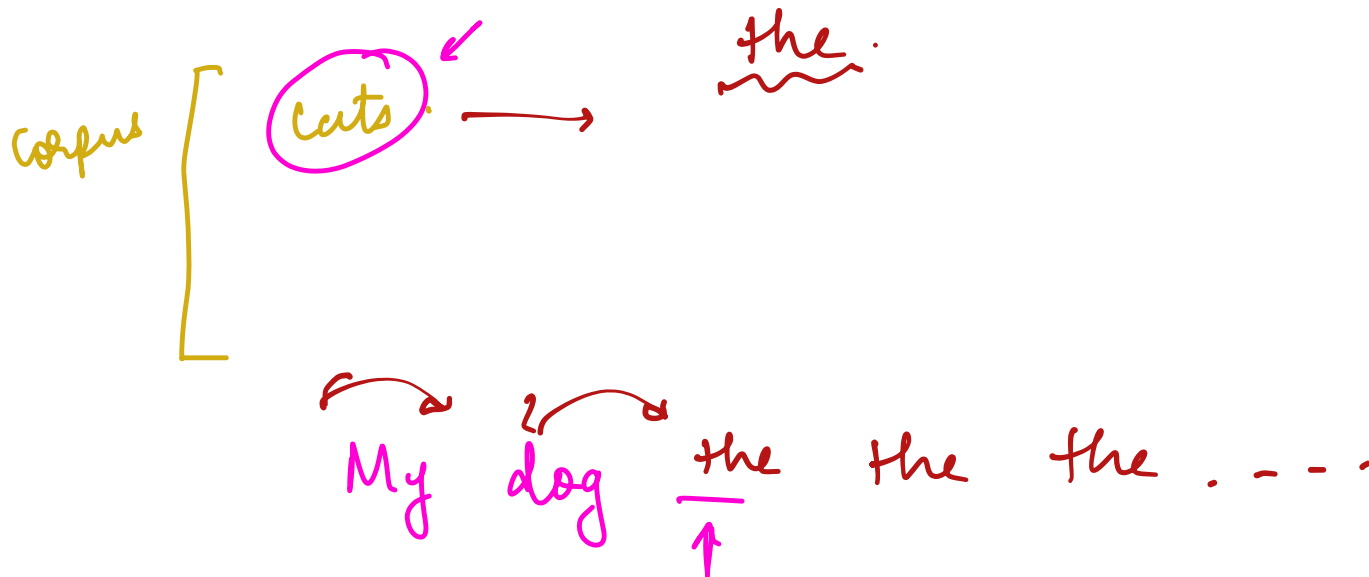
order?"

## My Pet.

My pets name is —

My pet is a —

My pet's color is ——

① $P(w_1, w_2, w_3, \dots w_n)$

$\rightarrow = P(w_1) \cdot P(w_2) \cdot P(w_3) \cdot P(w_4) \dots P(w_n).$

Naïve Bayes approach: All word occurrences are independent of the others.

corpus $\left[\begin{array}{l} \end{array}\right.$ cats $\longrightarrow$ the.

My dog the the the . - - -

① $P( \boxed{w_1 , w_2}^A , \underbrace{w_3}_{B} ) \xrightarrow{(a)} P(w_1) \cdot P(w_2) \cdot P(w_3)$

"Bigram"

$$\longrightarrow \underbrace{P(w_3 \mid \cancel{w_1} , w_2) \cdot P(w_1, w_2)}$$

$$\longrightarrow P(w_3 \mid w_2) \cdot P(w_2 \mid w_1) \cdot P(w_1).$$

$$\overset{w_1 \quad\quad w_2}{\overbrace{\underset{w_1 \quad\quad w_2}{\underbrace{\text{The man}}} \text{ } \underset{w_1 \quad w_2}{\underbrace{\text{saw}}} }} \text{ a girl looking through his}$$

telescope.

The $\boxed{\text{man}}$ ___

N-Gram → A window size of N is selected to calculate the probability.

# Bigram :

I have three books
with me.

$\langle start \rangle \rightarrow$ new token / new word we are
using to mark the start of
the sentence.

$P(w \mid \underline{\langle start \rangle})$

# Trigram :

$\rightarrow P(w \mid \langle start \rangle, \underline{\langle start \rangle})$

$\langle$ a    href $=$ " http : \\ . _ .    $\rangle$