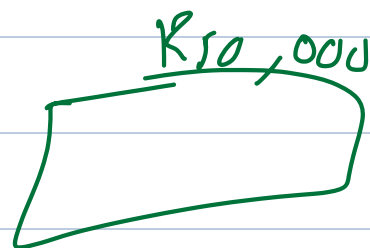
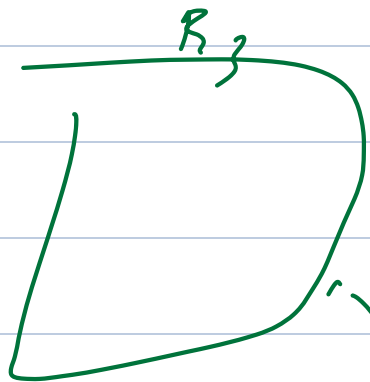
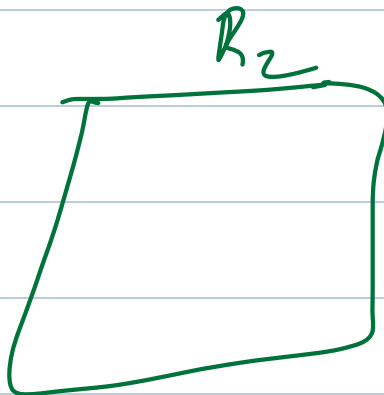
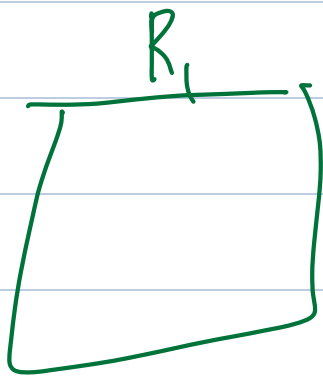


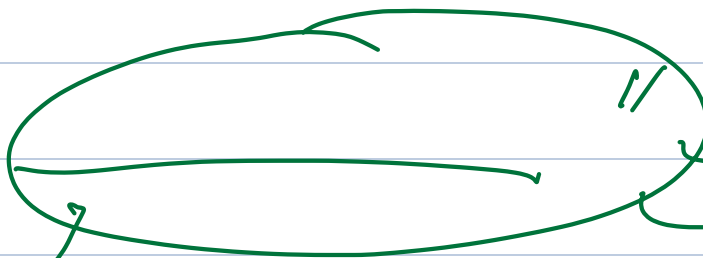
Language & topic modeling

Word2vec

skibgram



Query =



cosine (Query, R_1)

—, R_2)

—, R_3)

—, R_2)

abc

cde

def

↓
[---]

↓
[---]

↓
[---]

R1 \rightarrow 200 words

$$l = [0, 0, \dots, 0]$$

$l = l + \text{model}(\text{each word})$

$$l = l / 200$$



centroid

Query \rightarrow

\rightarrow origin of cosina

R1 $\cos(o, c) \quad \cos(o, c) \quad \cos(\text{cosina}, c)$

sum $\rightarrow 0.9$

centroid 1

centroid 2

centroid 10, 0, 0

R2 $\cos(o, c) + \cos(o, c) + \cos(\text{cosina}, c) \leftarrow 0.7$

King = Man + Prince

$\rightarrow _ + _ = _$

Language & Topic

Generative

Generative AI

Generative words

Training data CC0 corpus

please find my CV attached

→ $P(\text{attached} \mid \text{please find my CV})$

Context w_2

$P(w_1 \mid \text{context})$
 $P(w_2 \mid \text{context})$

$p(\text{Shakespeare} | \text{William})$

Why?

S1: the cat is small

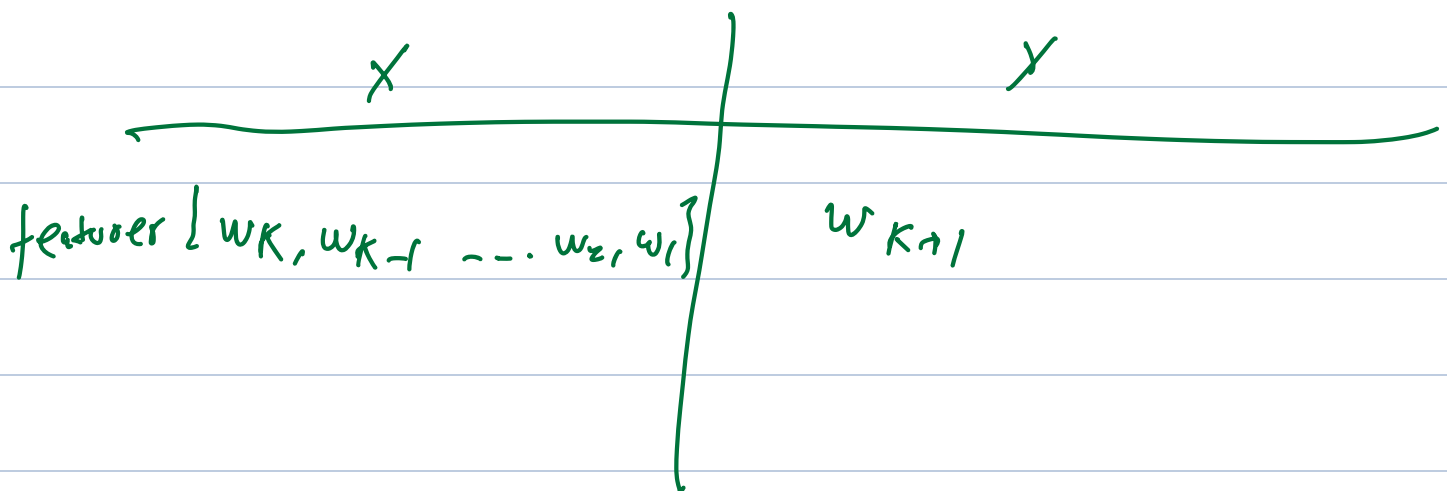
S2: smash the cat is

order of word matters

capital of france Paris

$p(\text{Paris} | \text{capital of france})$

$p(w_{k+1} | w_k, w_{k-1}, \dots, w_2, w_1)$



I

I hope

$P(\text{hope} | I)$

I hope that

$P(\text{that} | I \text{ hope})$

I hope that you

$P(\text{you} | I \text{ hope that})$

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

$$P(B|A) = \frac{P(A \cap B)}{P(A)}$$

'it's water is so transparent'

$$P('') = P(\text{it's}) \cdot P(\text{water} | \text{it's}) \cdot P(\text{it's} | \text{water}) - - -$$

$P(\text{transparent} | \text{it's water is so})$

$$= \frac{\text{count}(\text{it's water is so transparent})}{\text{count}(\text{it's water is so})}$$

Markov property

$P(\text{next word} | \text{ })$
 $\hookrightarrow \text{'K words'}$

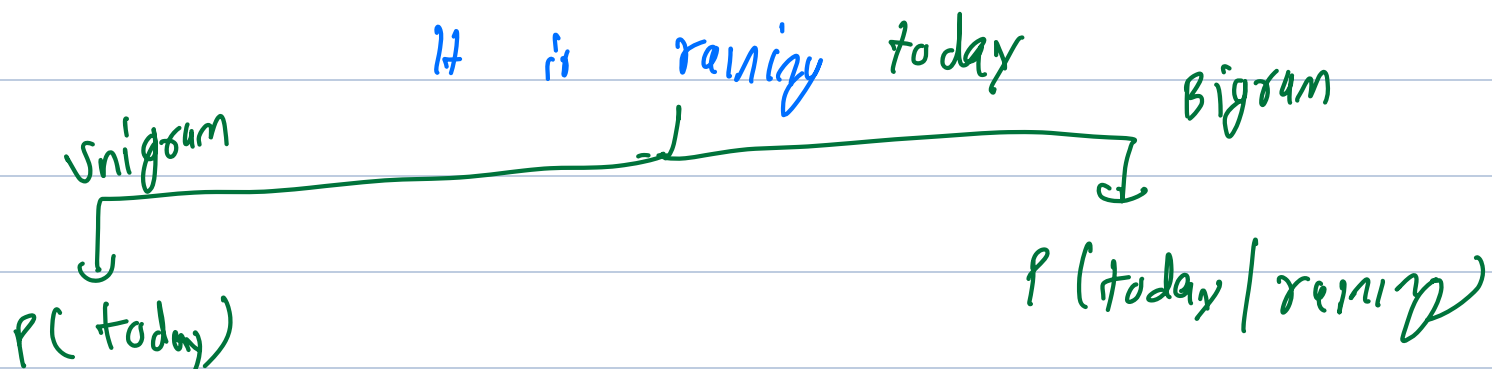
$$K=1 \quad P(\text{transparent} | \text{so})$$

$$K=2 \quad P(\text{transparent} | \text{is so})$$

$K=1 \Rightarrow$ Unigram

$K=2 \Rightarrow$ Bigram

$K=n-1 \Rightarrow$ N-gram model



$n = \text{history}$

Unigram

$$P(I \text{ have a dream}) \\ = P(I) \cdot P(\text{have}) \cdot P(a) \cdot P(\text{dream})$$

Bigram

$n=2$ <start>

$P(I \text{ have a dream})$

$$P(I | \boxed{}) * P(\text{have} | I) * P(a | \text{have}) * P(\text{dream} | a)$$

\downarrow
~~<start>~~

Corpus

<start> I want a dream job <start>

<start> I have a dog <start>

<start> I have a dream company <start>

of words
vocab size = 20

Unigram Table

<start>	I	want	a	dream	job	have	dog/country
3/20	3/10	1/20	-	-	-	-	-

$$P(\text{I have a dream}) = P(I) \cdot P(\text{have}) \cdot P(a) \cdot P(\text{dream})$$

$$= \text{---} \cdot \text{---} \cdot \text{---} \cdot \text{---}$$

Bigram Model

$$P(w_2 | w_1) \rightarrow P(\text{like} | I)$$

I like

$$P(\text{AI} | \text{gen})$$

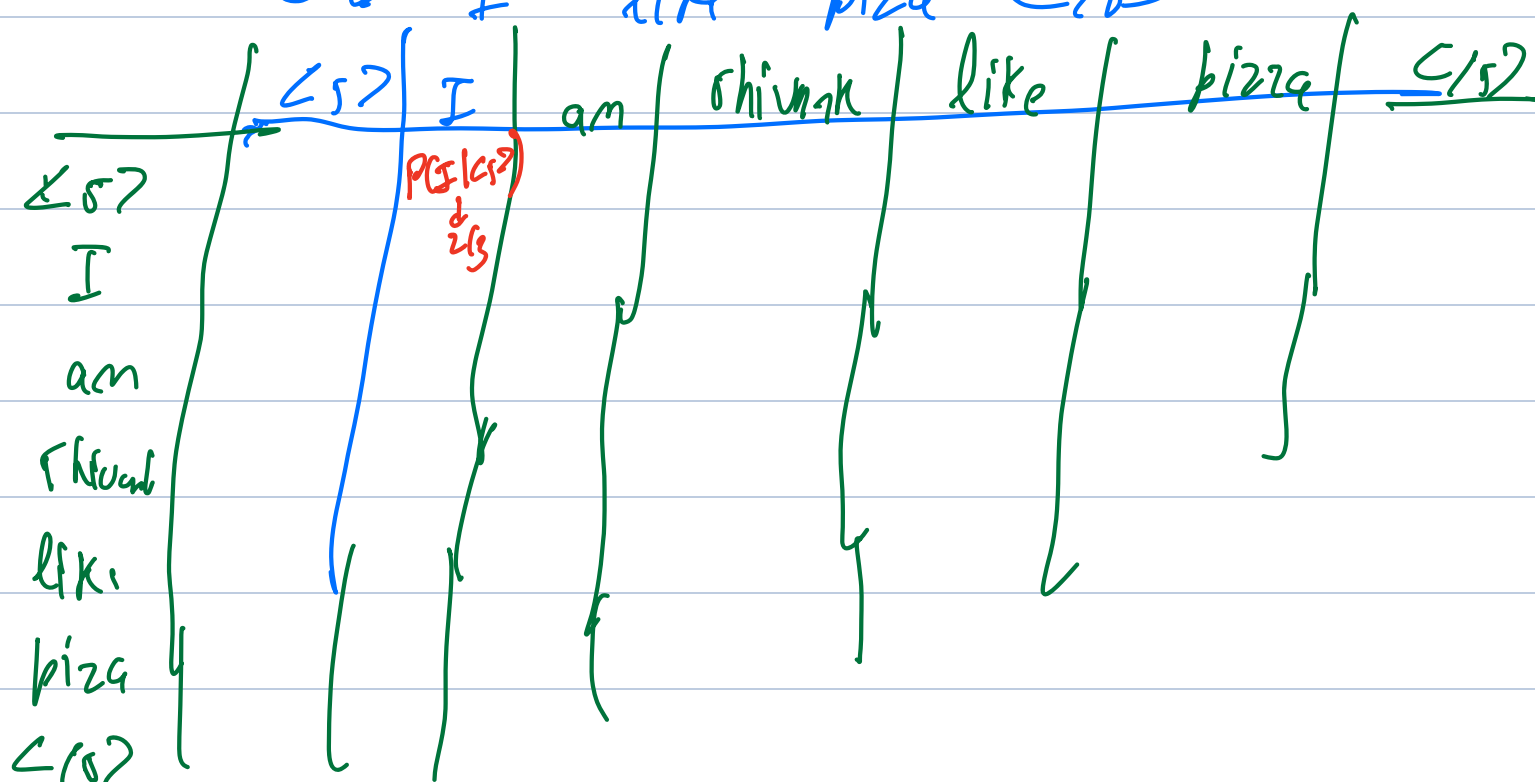
Generative AI

Coopus

<5> I am shivank </5>

<5> shivank I am </5>

<5> I like pizza </5>



Query \Rightarrow <5> shivank I like pizza </5>

$P(\text{shivank} | \langle 5 \rangle) \cdot P(I | \text{shivank}) \dots$

<r> I am shivank </r>

<r> shivank I am </r>

<r> I like pizza </r>



New Delhi is capital of India

I am in Mumbai

(is New) I

is Delhi I

is capital I

is of I

is Mumbai O

is in O