

Class

Recap

Reminders

word2vec

GloVe

ELMo, BERT, GPT

Recap

(x_1, \dots, x_n) input sequence

(y_1, \dots, y_n) output sequence

- associate inputs with outputs

I love you a lot

Te amigo mucho

- variable length
- long range dependency

$$q_i = W_q x_i$$

$$k_i = W_k x_i \quad \text{for } i \in [n]$$

$$v_i = W_v x_i \quad = \{1, 2, \dots, n\}$$

$$w_{ij} = q_i^T k_j$$

$$\text{Softmax}(w_{i1}, w_{i2}, w_{i3}, \dots, w_{in})$$

$$\frac{\exp(w_{ij})}{\sum_{j'} \exp(w_{ij'})}$$

$$y_i = \sum_{j=1}^n w_{ij} v_j$$

$$y_i \in \mathbb{R}^d$$

$$v_j \in \mathbb{R}^d$$

$$q_i \in \mathbb{R}^{q\text{-dim}}$$

$$x_i \in \mathbb{R}^{in\text{-dim}}$$

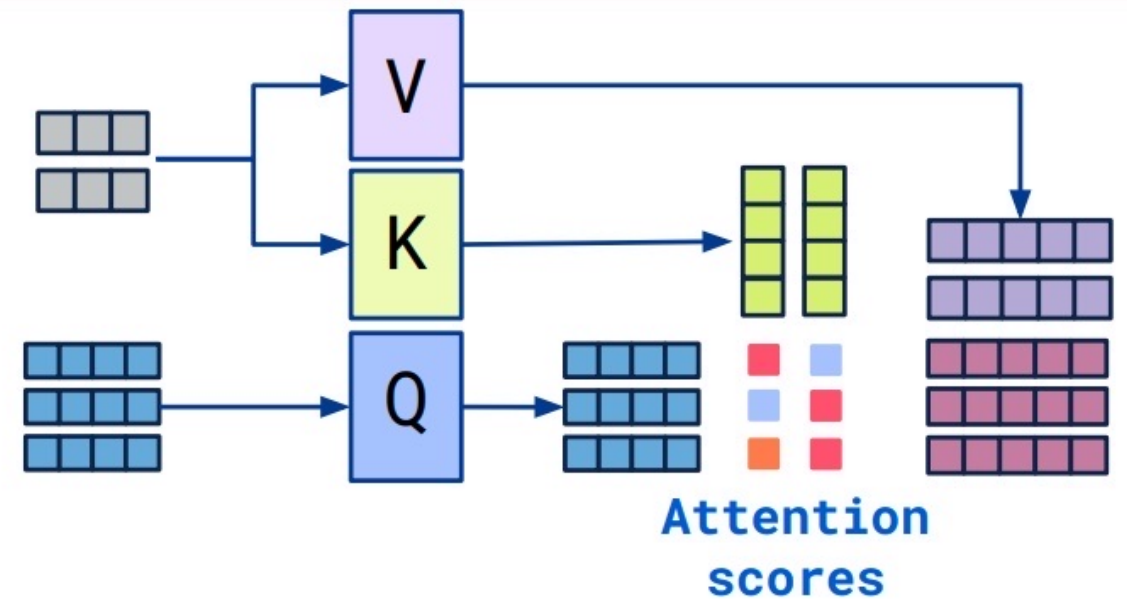
$$W_q, W_k$$

$$q\text{-dim} \times in\text{-dim}$$

$$W_v$$

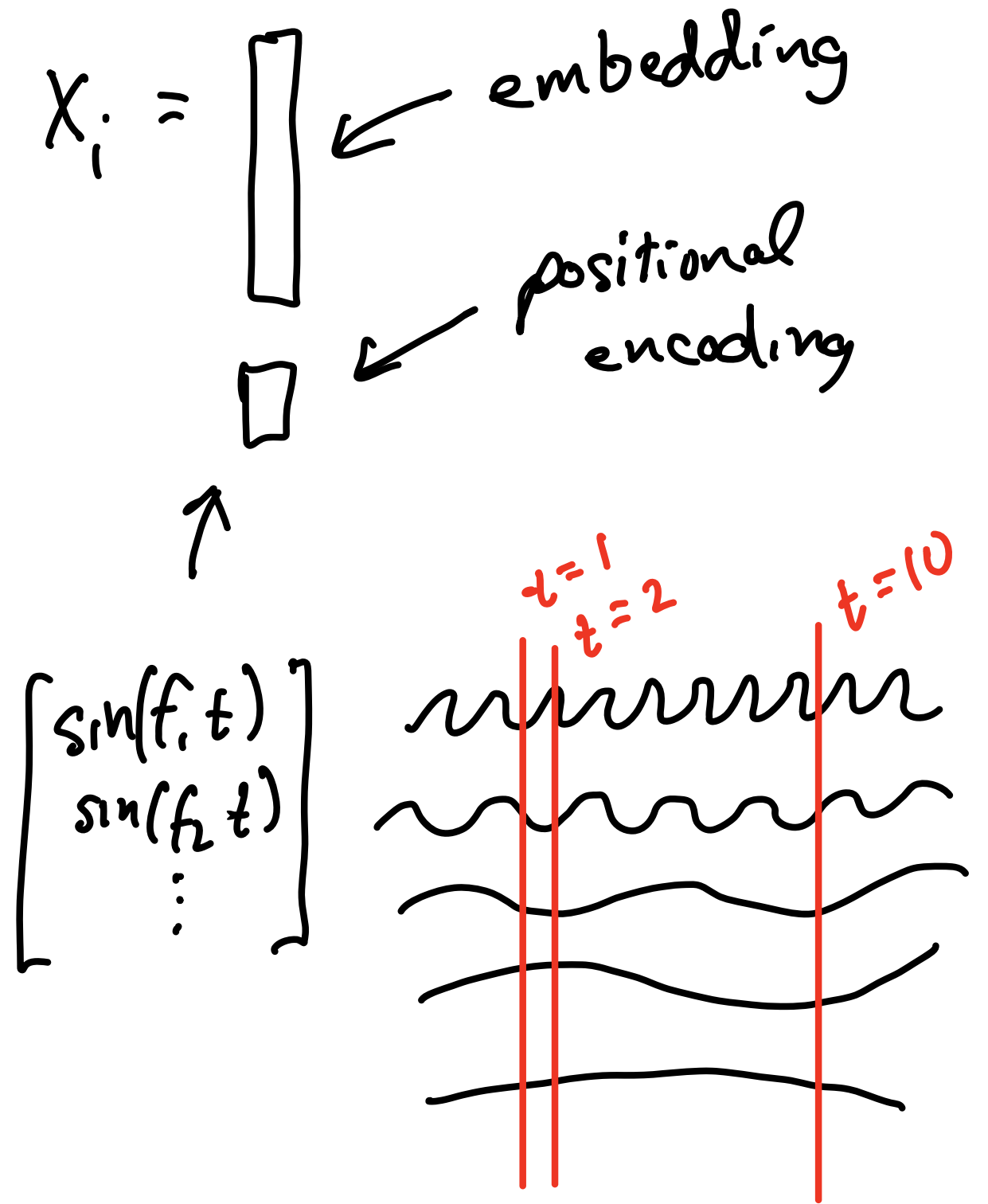
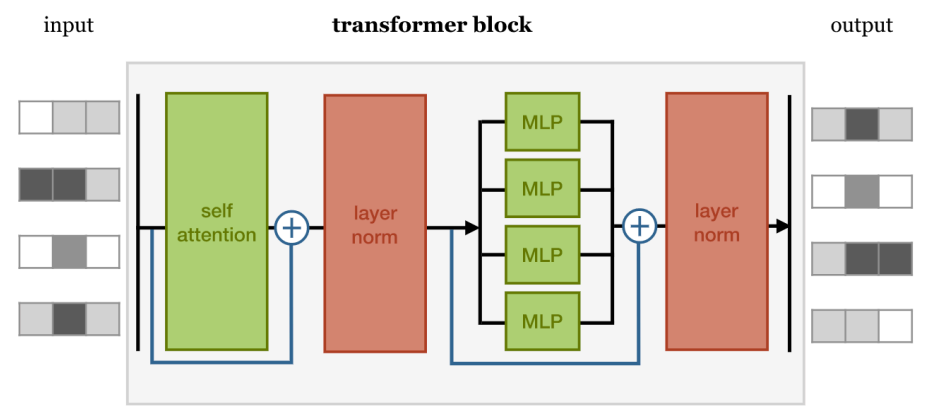
$$d \times in\text{-dim}$$

(x_1, \dots, x_n)
 (z_1, \dots, z_d)



Cross attention

Transformer



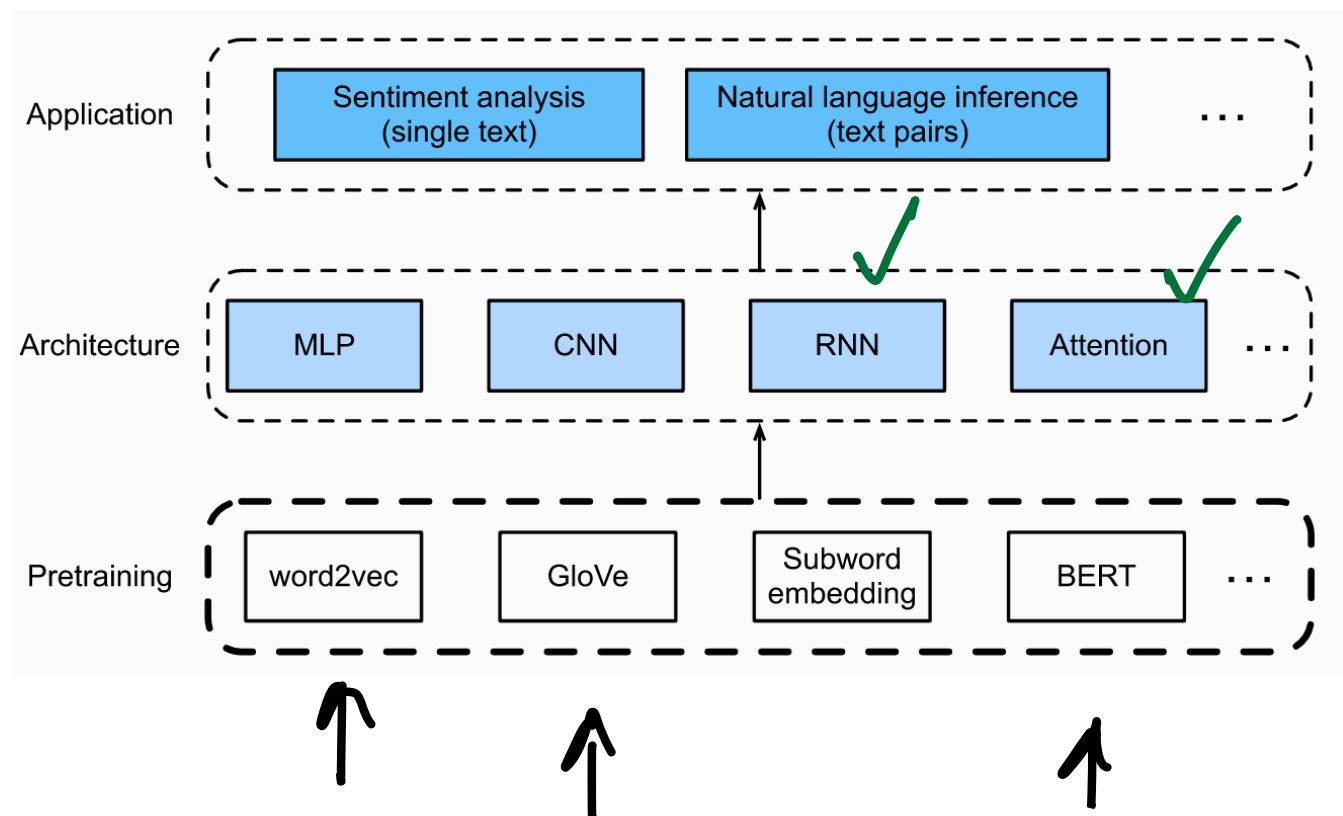
$$\|a\|_2 = \beta \iff \|a\|_2^2 = \beta^2$$

$$\|a\|_2^2 = a^T a = \beta^2$$

Reminders

- Form 23/26
- Homework due Spm Friday
↳ self grade due Monday
- Proposal due Spm Monday

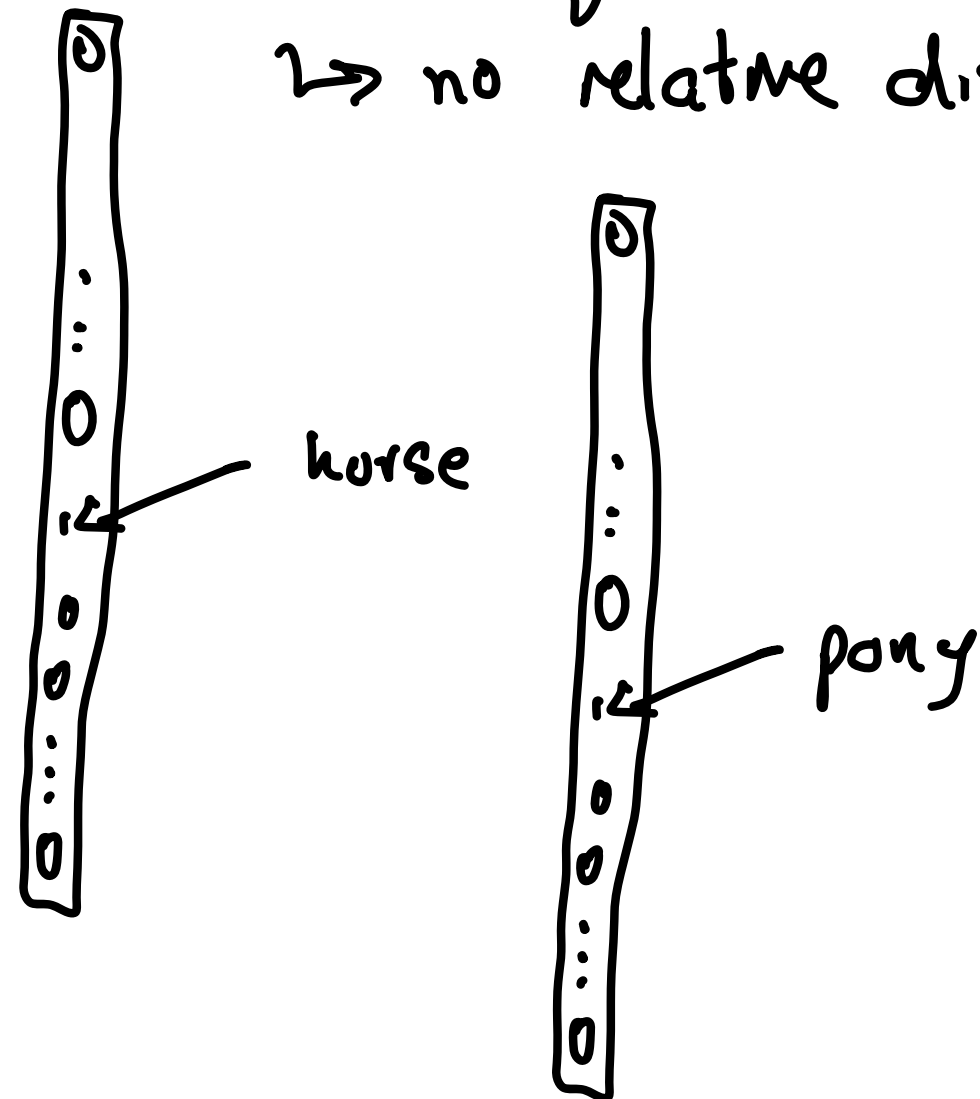
Overview



word \rightarrow vector
(meaningful)

One Hot Encoding

\hookrightarrow many dimensions
 \hookrightarrow no relative distances



word2vec

2-gram

Gatsby believed in the green light, the orgastic future that year by year recedes before us. It eluded us then, but that's no matter—tomorrow we will run faster, stretch out our arms farther. . . . And one fine morning—

So we beat on, boats against the current, borne back ceaselessly into the past.

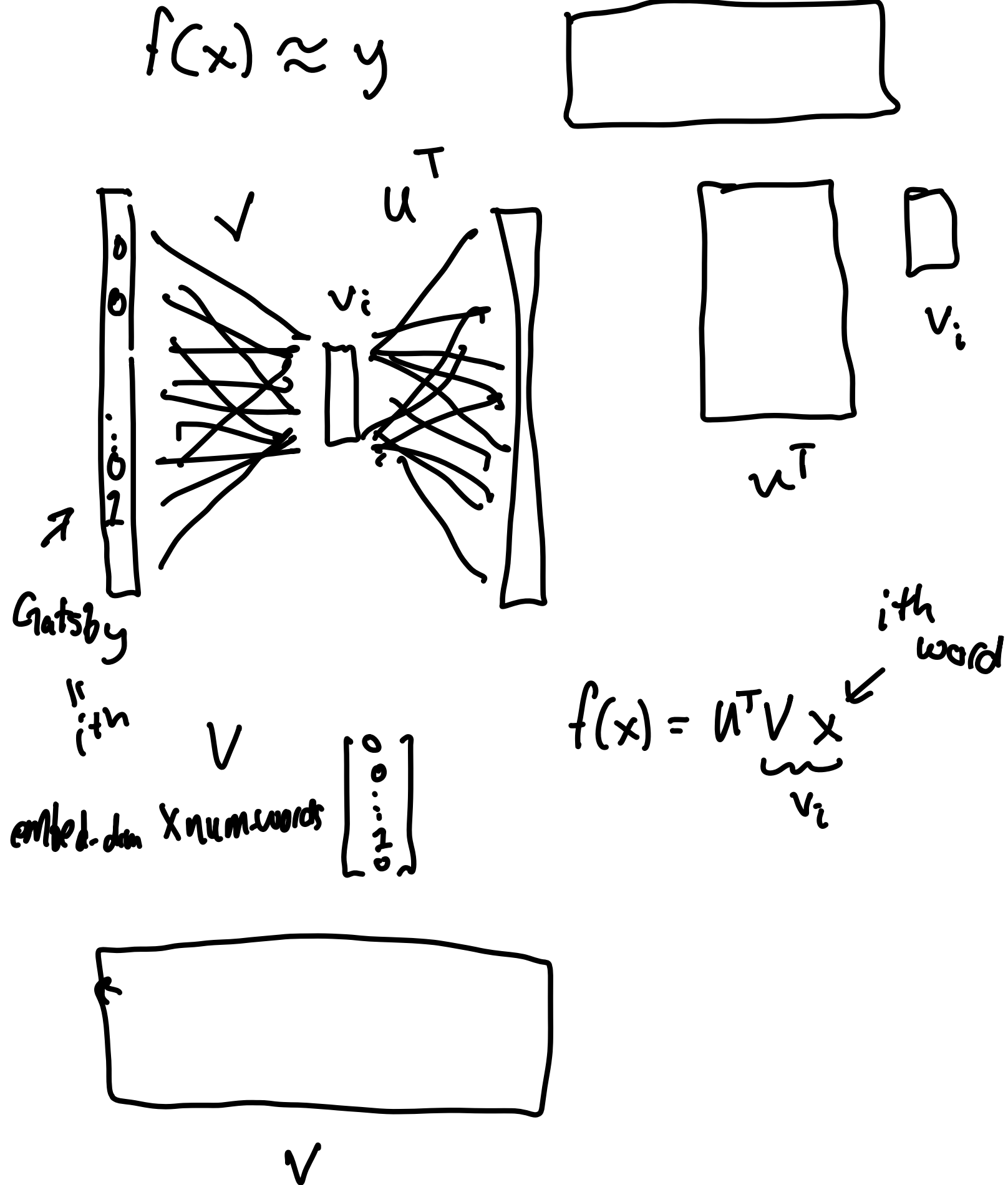
problems

↳ too sparse

↳ throw away half the model

- x, y
- (Gatsby, believed)
- (believed, in)
- (in, the)
- (the, green)
- ⋮

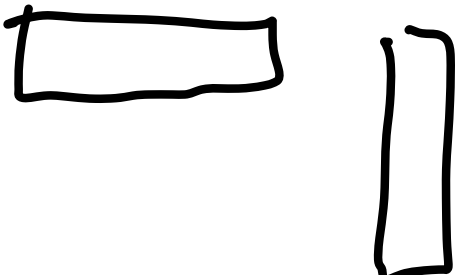
(Gatsby, in)



$f(x)$ ← embedding of x

(x, y) $f(x) \approx f(y)$

$f(x)^T f(y)$ large
if x is
close in
meaning
to y



$$f(x) = Vx$$

$$\mathcal{L}(v) = -\sum_{x,y} f(x)^T f(y)$$

↑ positive pair

$$+ \sum_{x,z} (f(x)^T f(z))^2$$

↑
negative
pairs

Extension: (x, y)

also include neighbors of
neighbors

n -gram vs z -gram

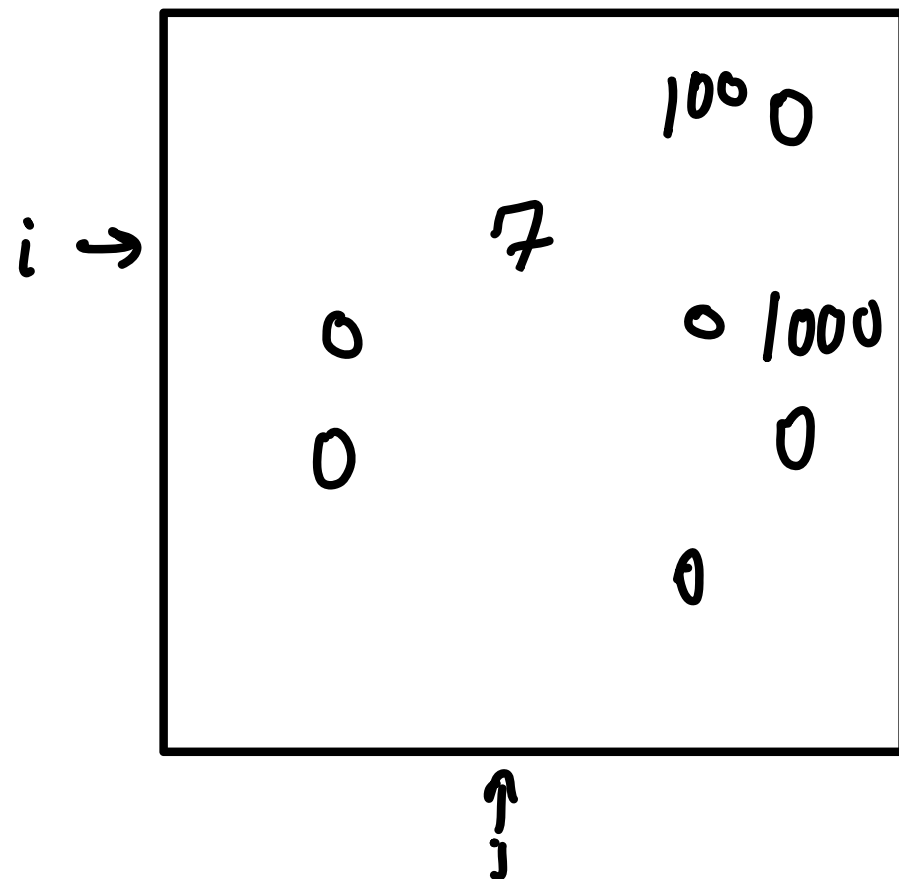
Glove

Global Vectors

X = co-occurrence matrix

n = # unique words

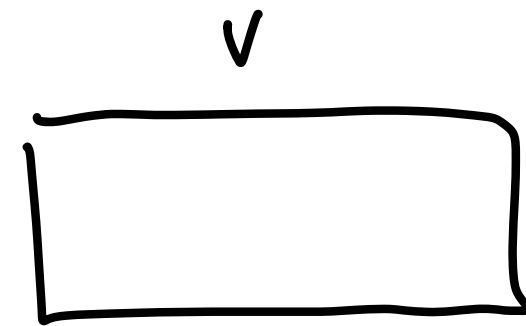
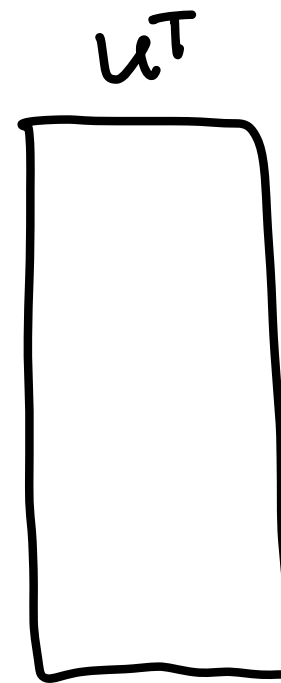
$X \in \mathbb{R}^{n \times n}$



$$\log(X) \approx U^T V$$

$n \times h \quad h \times n$

trainable parameters
 $U, V \in \mathbb{R}^{h \times n}$
 $c, b \in \mathbb{R}^n$



$$\log(x_{ij}) \approx u_i^T \cdot v_j$$

embedding v_i for word i

$= v \times$ ← one hot encoded vector for word i

$$\mathcal{L} = \frac{1}{2} \sum_{i,j} (u_i^T v_j + b_i + c_j - \log(x_{ij}))^2 f(x_{ij})$$

really want to only look at non zero

↑
embedding

close to 1 for large values, close to 0 for small values

Embedding with context

↳ crane

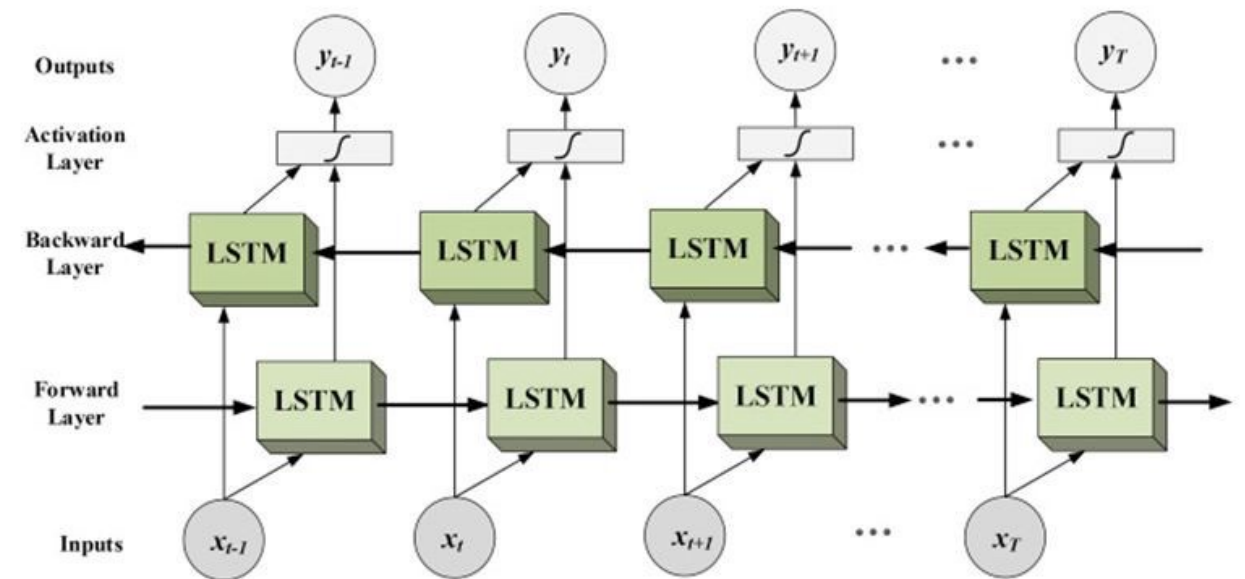
↳ date

↳ dough

ELMo

Embeddings from Language Models

Bidirectional LSTM



Idea: Replace linear layer in word2vec with fancier architecture

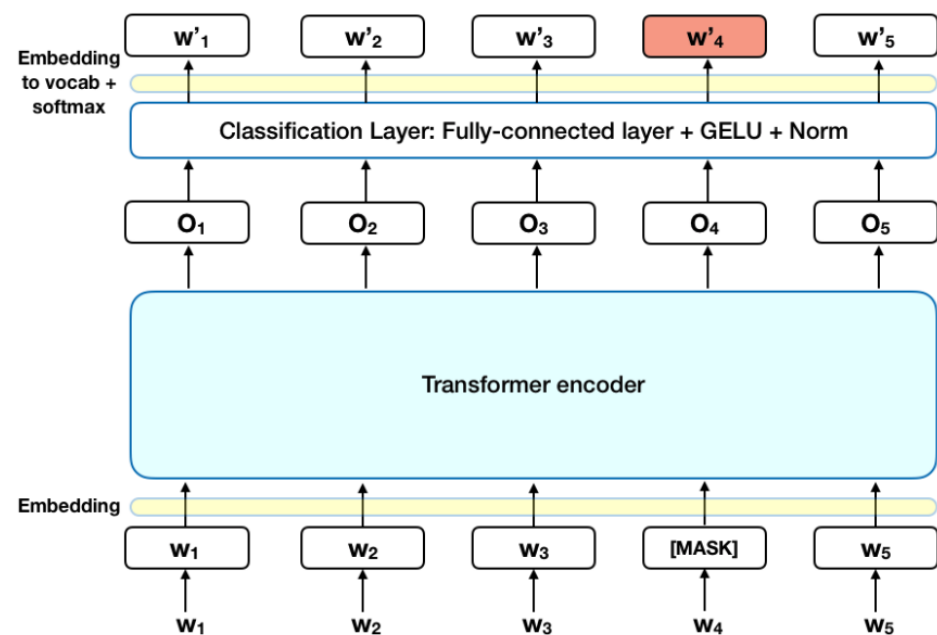
BERT

Bidirectional

Encoder

Representations from
Transformers

Idea: Replace with
transformer



↳ still for next sentence

↳ drop out random words

↳ piece to kenitization

walking → walk,ing

GPT

↳ masked tokens

↳ huge data set and
lots of parameters