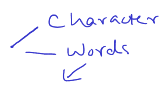
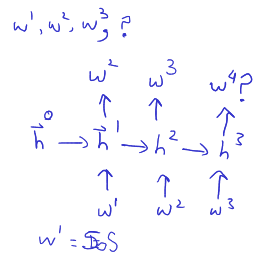


Language Modelling

① Text Prediction



Given words w^1, w^2, \dots, w^t , predict next word w^{t+1} .



Input

- Consider the English language. Let English vocabulary V have K words.
- Each word can be represented by a K -dimensional vector with 1-of- K coding.
 - ↳ Tremendously inefficient. But how else can we convert a word into a vector of numbers that we can input into a neural network?
 - ↳ Each word vector is orthogonal to every other word vector. This is not realistic at all since words have similarities.

Workaround: Project word vectors into a lower-dimensional space via a projection/embedding matrix E of size $D \times K$ where $D < K$.

Matrix E can be learned.

Output

- Output layer = softmax on K neurons.
 - $y_k^{t+1} = p(V_k | w^1, w^2, \dots, w^t)$
 - = prob. of k -th word at time $t+1$ given previous words
- $\vec{y}^{t+1} = \begin{bmatrix} p(V_1 | w^1, \dots, w^t) \\ p(V_2 | w^1, \dots, w^t) \\ \vdots \\ p(V_K | w^1, \dots, w^t) \end{bmatrix}$
 ↑
 prob. distribution $K \times 1$

Loss

Cross-entropy between output words and target words.

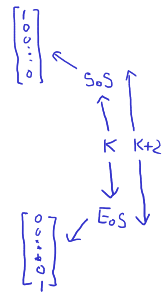
$$L(\underbrace{(\vec{d}^1, \vec{d}^2, \vec{d}^3, \dots, \vec{d}^T)}_{\text{1-of-K coded vectors}}, \underbrace{(\vec{y}^1, \vec{y}^2, \vec{y}^3, \dots, \vec{y}^T)}_{\text{prob. distribution vectors}}) = - \sum_{t=1}^T \sum_{k=1}^K d_k^t \ln y_k^t = - \sum_{t=1}^T \ln y_{\text{index corresponding to } w^{t+1}}^t$$

d_k^t 1 for desired word at time t
 0 for the rest.

Training Prop using millions of sentences. Each sentence is 1 training sample.

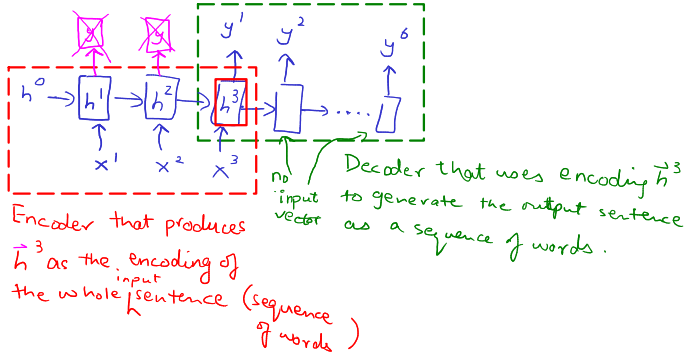
Sentence Language Generation: Next word w^{t+1} is drawn from prob. dist. \vec{y}^t

- Feed w^{t+1} as next input and find \vec{y}^{t+1} .
 - Repeat until ending symbol is drawn.
- 1-of-($K+2$) coding.

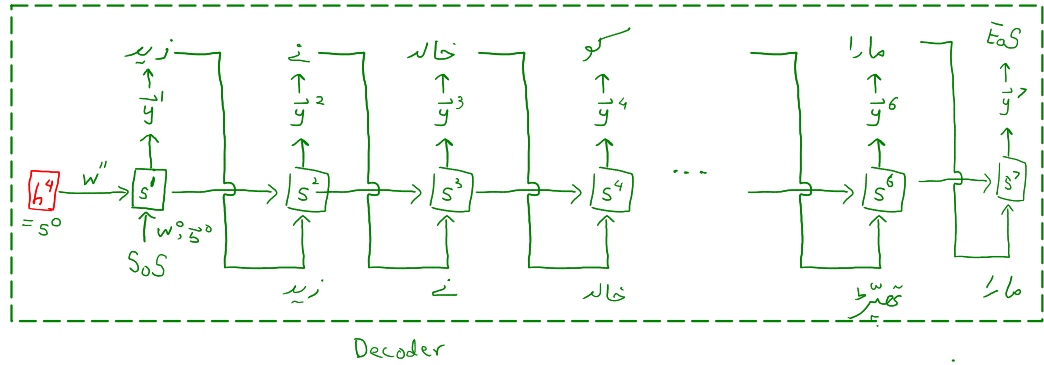
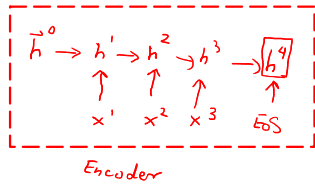


Language Translation

SoS Zaid slapped Khalid EoS \rightarrow EoS زید نے خالد کو تھپڑ مارا EoS
 $x^0 \ x^1 \ x^2 \ x^3 \ x^4 \ \rightarrow \ y^7 \ y^6 \ y^5 \ y^4 \ y^3 \ y^2 \ y^1 \ y^0$



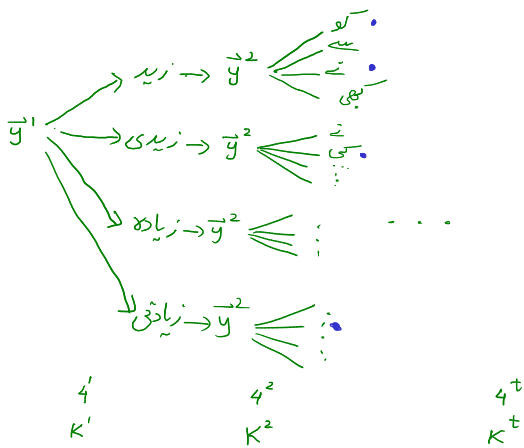
A better decoder: make prob. dist \vec{y}^{t+1} depend on word drawn from \vec{y}^t as well.



Draw a sample from \vec{y}^t .

- ① Argmax \vec{y}^t
- ② Top-K (\vec{y}^t).

$$y_j^t = P(O^t = V_j \mid \underbrace{O^{t-1}, O^{t-2}, \dots, O^1}_{\text{all words output so far}}, \underbrace{I^1, I^2, \dots, I^N}_{\text{all input words}})$$



Beam Search

- Terminate each branch when it produces EoS word.
- We let all Top-K branches continue until termination.
- From all the terminated branches, pick the one with the highest overall probability.

$$\text{Top-K } P(O^2 O^1 \mid I^1, I^2, \dots, I^N) = \underbrace{P(O^2 \mid O^1, I^1, I^2, \dots, I^N)}_{y^2 \text{ index}} \underbrace{P(O^1 \mid I^1, I^2, \dots, I^N)}_{y^1 \text{ index}}$$

Product rule: $P(AB|C) = P(A|BC)P(B|C)$

At time 2, this will give me overall probabilities of 16 branches.
 From those 16 branches pick the top-K.