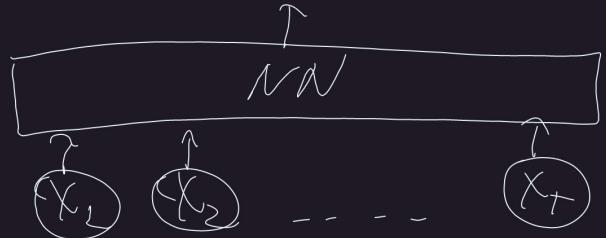


NLP for texts



Lemmatization

"went" → "go"
"better" → "good"
"corpus" → "corpus"

Stemming

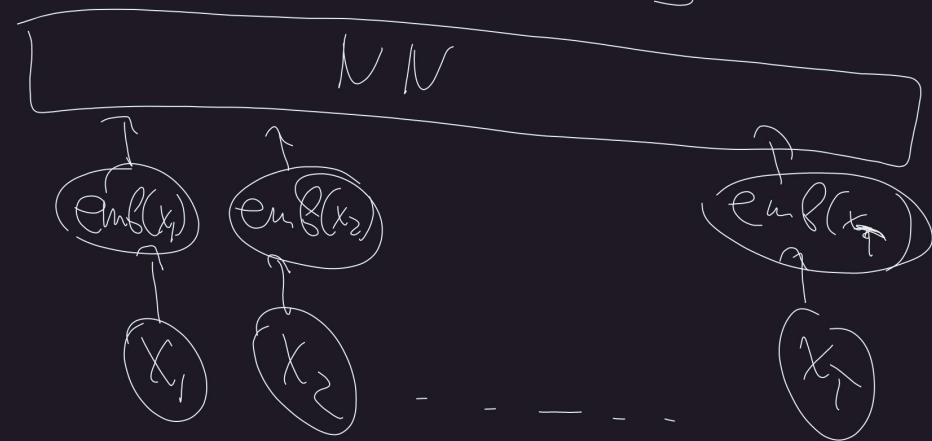
models
modelling
modelled
- - -
stemming model

$\text{lem}(\text{"modelling"}, \text{pos} = \text{"V"}) \approx$
 $\text{lem}(\text{"modelling"}, \text{pos} = \text{"S"}) = \text{"model"}$

stop words: "the", "a", "an", "and", "in"

pipeline: words \rightarrow stemming \rightarrow remove least
lemmatization freq. tokens \rightarrow remove stop words

$$x_t \in \{1, 2, \dots, V\} \rightarrow \text{emb}(x_t) \in \mathbb{R}^{\text{emb-dim}}$$



<UNK>, <PAD>, ...

(Byte pair encoding)
BPE tokeniser

Corpus

- "hug" - 10
- "pug" - 5
- "fun" - 12
- "fun" - 4
- "hugs" - 5

Vocab

- "h", "u", "g", "
- "p", "n", "f", "
- "s", "

All pairs

- "hu" - 15
- "hg" - 0
- "ug" - 20
- ..
- ..
- ..

difference:

"hug" → ("h", "u", "g")
 ↓
 "h", "g", "s" → "h", "g", "s", "

Vocab.

- "h", "u", "g", "
- "p", "n", "f", "
- "s", "

("gs"), ("hug")

Word Piece Tokenizer

"word" → "w", "#to", "#z", "#d"

$$\text{Score}(\text{token}_1, \text{token}_2) = \frac{\text{freq}(\text{token}_1, \text{token}_2)}{\text{freq}(\text{token}_1) \cdot \text{freq}(\text{token}_2)}$$

"Un", "#able" → "#bo", "#z" → "#or"
"fre", "#men" → "fremea", "w" | "#o" → "wo"

n-gram tokenizer

"hug" - 10 → p_1
"hug" - 5 → p_2
"hug" - 4 → p_3
"hug" - 3 → p_4

"hug" → p_1 , "u" → p_2 , "g" → p_3 , "h" → p_4

$$\text{P}(\text{"hug"}) = \text{P}(\text{"h"}, \text{"u"}, \text{"g"}) = \text{P}(\text{"h"}) \cdot \text{P}(\text{"u"}) \cdot \text{P}(\text{"g"})$$

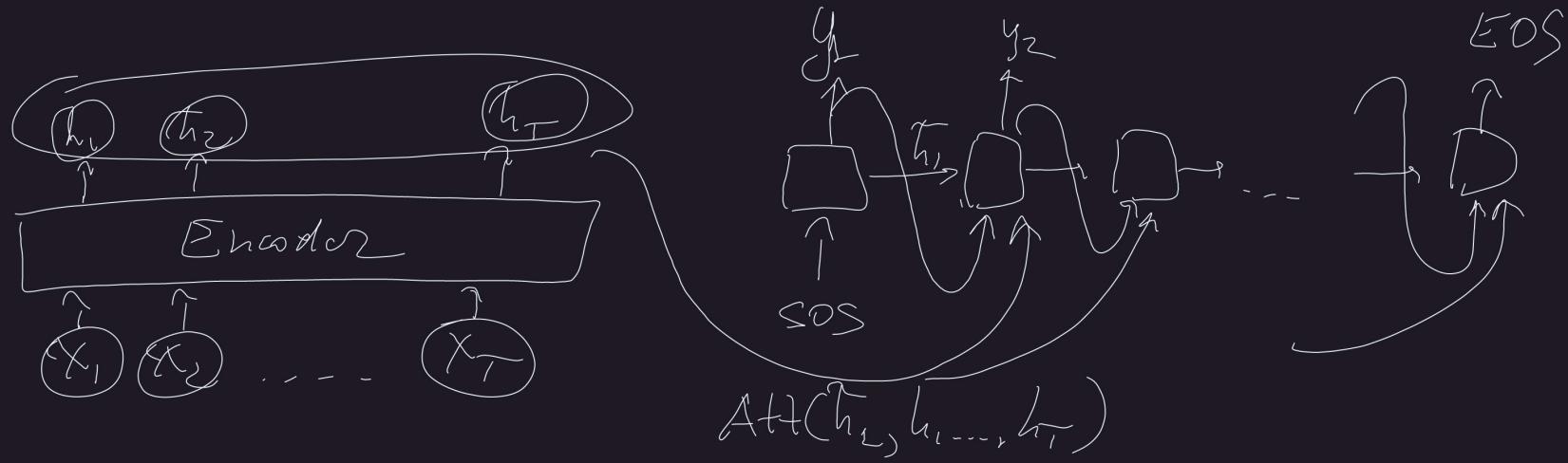
"hugs" - 5 → p_1

$$\text{loss} = -(\log p_1)(0) - (\log p_2)5 \dots$$

$$= \frac{15}{210} - \frac{20}{210}$$

inference:

"hugs" → "h", "u", "g", "s"



teacher forcing
exposure bias fix

$$\begin{aligned}
 p(y_1, \dots, y_T | x, \theta) &= \prod_t p(y_t | x, \theta) = p(y_1 | x, \theta) \cdot p(y_2 | y_1, x, \theta) \cdots \\
 &= \prod_{t=1}^T p(y_t | y_1, \dots, y_{t-1}, x, \theta)
 \end{aligned}$$

$$\log p(y_1 \dots y_T | x, \theta) + \lambda \log p_M(y_1 \dots y_T | x, \theta) + \beta t$$

beam search

$$\arg \max_y p(y | x, \theta) = \arg \max_y \prod_{t=1}^T p(y_t | y_1 \dots y_{t-1}, x, \theta)$$

$$\text{top-2} \left(\begin{array}{l} "a" - 0.5 \\ "b" - 0.4 \\ "c" - 0.1 \end{array} \right)$$

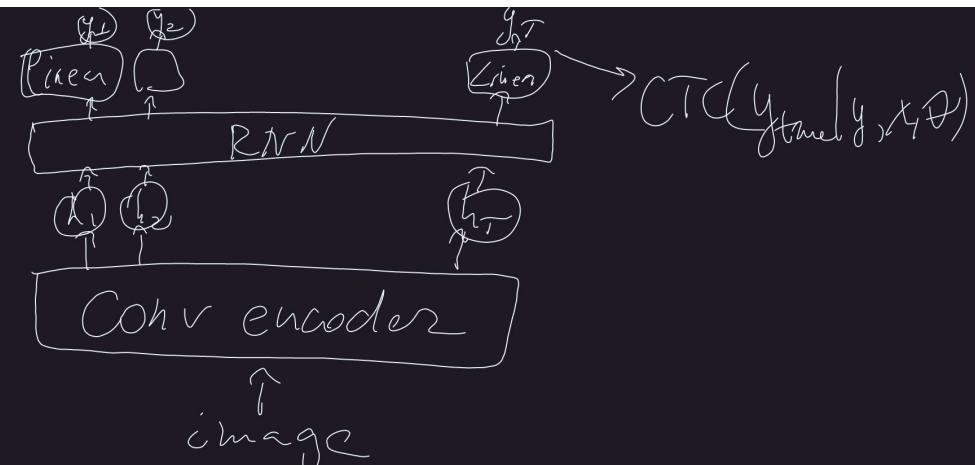
$$\begin{array}{ll} \text{t=1} & \text{t=2} \\ \begin{array}{l} "aa" \\ "ab" \\ "ac" \end{array} & \begin{array}{l} 0.5, 0.5 \\ 0.5, 0.4 \\ 0.5, 0.1 \end{array} \end{array} \text{top-2}$$

$$\begin{array}{ll} T, V & O(\sqrt{T}) \\ \text{t=3} & \\ \begin{array}{l} "aaa" \\ "aab" \\ "aac" \\ "aba" \\ "abb" \\ "abc" \end{array} & \begin{array}{l} "aaa" \\ "aab" \\ "aac" \\ "aba" \\ "abb" \\ "abc" \end{array} \end{array}$$

K-beam size $O(T \cdot K \cdot V)$

CTC loss

handwritten text rec.



[hello] → "h", "e", "e", "l", "l", "o"

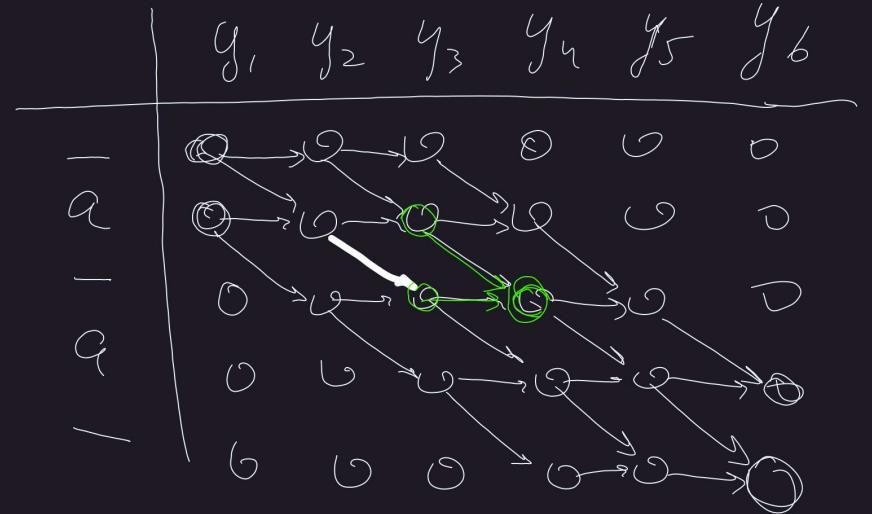
CTC-loss:

$$p(y_{\text{true}} | x, \theta) = \sum_{y: \text{enc}(y) = y_{\text{true}}} \prod_{t=1}^T p_t(y_t | x, \theta)$$

CTC(- - h - e e - l l o -) =

= h - e - l - l - o - → hello

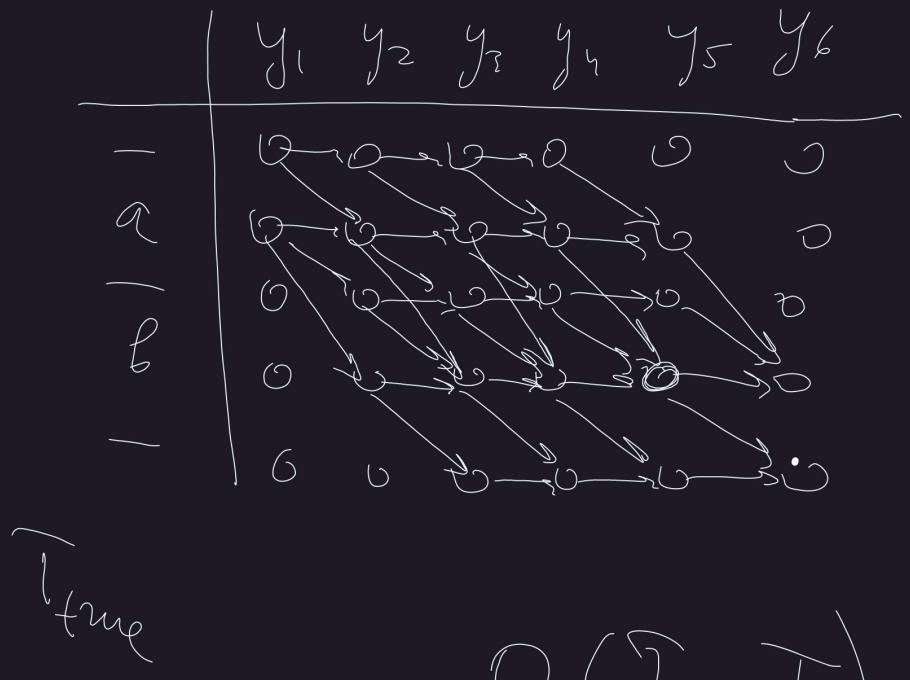
$y_{true} = "aa"$, $T=6$



$$TC(\underline{\dots} _ a _ a) = aa$$

$$d_{s,t} = p(y_t | x_t, \theta) (d_{s,t+1} + d_{s-1,t+1})$$

$$y_{\text{true}} = "ab", T=6$$



$$\Omega \left(\lceil \frac{n}{\epsilon} \rceil T \right)$$

$$\mathcal{L}_{s,t} = \left(\mathcal{L}_{s,t-1} + \mathcal{L}_{s-1,t-1} + \mathcal{L}_{s-2,t-1} \right) p_t(y_t | x_t, \theta)$$

if $y_{\text{true},s} \neq y_{\text{true},s-2}$
and $y_{\text{true},s} \neq -$

$$\mathcal{L}_{s,t} = \left(\mathcal{L}_{s,t-1} + \mathcal{L}_{s-1,t-1} \right) p_t(y_t | x_t, \theta)$$

otherwise

$$CTC = \mathcal{L}_{4,6} + \mathcal{L}_{5,6}$$

Non-diff. losses

$\text{BLEU}(y, y_{\text{true}})$ — machine translation

$\text{CIDEr}(y, y_{\text{true}})$ — image2caption

$\sum_i p(x_{y(i)}, x_{y(i+1)})$ — TSP



$$F(\theta) = \underset{p(y|x, \theta)}{\mathbb{E}} f(y) \xrightarrow{\text{non-dif.}} \min_{\theta}$$

$$\begin{aligned} \nabla_{\theta} F(\theta) &= \nabla_{\theta} \underset{p(y|x, \theta)}{\mathbb{E}} f(y) = \sum_{\theta} \sum_y p(y|x, \theta) f(y) = \\ &= \sum_y \nabla_{\theta} p(y|x, \theta) \cdot f(y) = \sum_y p(y|x, \theta) \cdot \frac{\nabla_{\theta} p(y|x, \theta)}{p(y|x, \theta)} f(y) = \\ &= \underset{p(y|x, \theta)}{\mathbb{E}} \nabla_{\theta} \log p(y|x, \theta) f(y) \approx \\ &\approx \frac{1}{M} \sum_{j=1}^M \nabla_{\theta} \log p(y_j|x, \theta) f(y_j), \quad y_j \stackrel{\text{i.i.d.}}{\sim} p(y|x, \theta) \end{aligned}$$

baseline $b(x)$

$$\nabla_{\theta} F(\theta) = E_{p(y|x, \theta)} \nabla_{\theta} \log p(y|x, \theta) (f(y) - b(x))$$

$$\begin{aligned} E_{p(y|x, \theta)} \nabla_{\theta} \log p(y|x, \theta) b(x) &= \sum_y p(y|x, \theta) \frac{\nabla_{\theta} p(y|x, \theta)}{p(y|x, \theta)} b(x) \\ &= b(x) \sum_y (\nabla_{\theta} p(y|x, \theta)) = b(x) \nabla_{\theta} \left(\sum_y p(y|x, \theta) \right) = 0 \end{aligned}$$

REINFORCE / log-derivative trick

Good choices for $b(x)$:

$$① \quad b(x) = E_{p(y|x, \theta)} f(y)$$

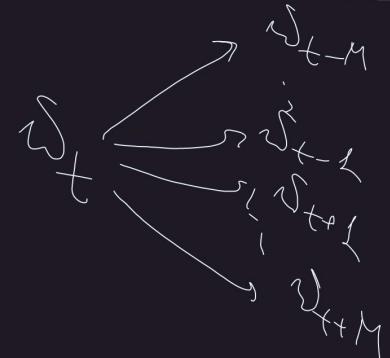
$$② \quad b(x) = f(y_{\text{greedy rollout}})$$

Word Embedding

$x_t \in \{1, 2, \dots, V\} \rightarrow \text{emb}(x_t) \in \mathbb{R}^{\text{emb_dim}}$

Word2Vec

$$\left[\begin{matrix} \omega_1 & \omega_2 & \omega_3 & \dots \end{matrix} \right]$$



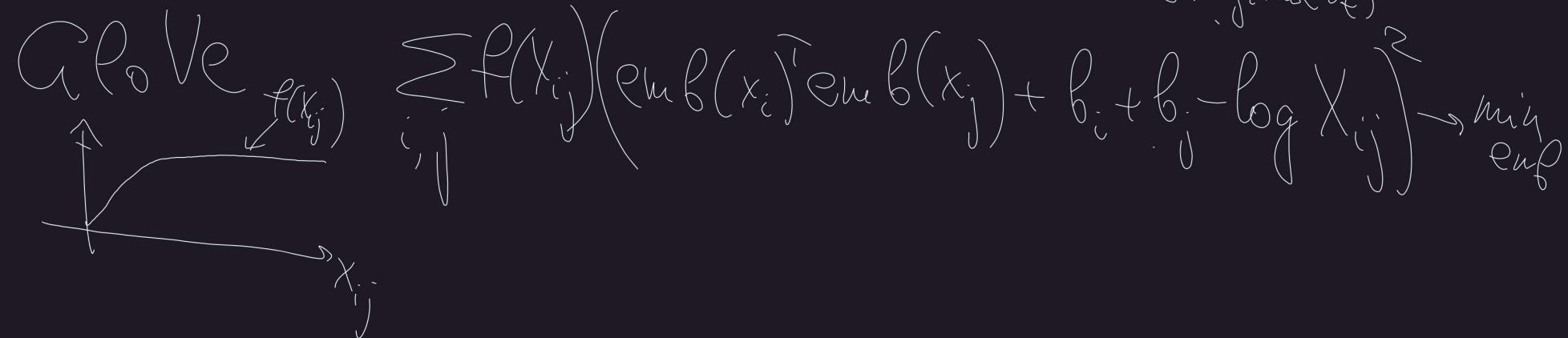
$$\exp(\text{emb}(\omega_c)^T \text{emb}(\omega_t)) \\ \sum_{j=1}^V \exp(\text{emb}(\omega_j)^T \text{emb}(\omega_t))$$

$$\sum_t \sum_{c \in C} \log p(\omega_c | \omega_t) \rightarrow \max_{C \in Q}$$

negative sampling

$$\sum_t \left[\sum_{c \in C} \log \sigma(\text{emb}(\omega_c)^\top \text{emb}(\omega_t)) + \sum_{n \in N_{t,c}} \log \sigma(-\text{emb}(\omega_n)^\top \text{emb}(\omega_t)) \right] \rightarrow \max_{\text{emb}}$$

fast-text : score(\(\omega_c, \omega_t\)) = \(\text{emb}(\omega_c)^\top \sum_{w \in \text{neighbors}(\omega_t)} \text{emb}(w)\)



Cove

ELMo

BERT

