Sequence to sequence models

Beam search

deeplearning.ai

# Beam search algorithm

$B = 3$   (beam width)

## Step 1

$$\rightarrow P(y^{<1>} \mid x)$$

$$
\begin{bmatrix}
\text{a} \\
\vdots \\
\text{in} \\
\vdots \\
\text{jane} \\
\vdots \\
\text{september} \\
\vdots \\
\text{zulu}
\end{bmatrix}
$$

10000



$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow \square \rightarrow \hat{y}^{<1>}$

$x^{<1>}$   $x^{<T_x>}$

Andrew Ng

# Beam search algorithm

$(B = 3)$

## Step 1    Step 2

$\begin{bmatrix} a \\ \vdots \\ in \\ \vdots \\ jane \\ \vdots \\ september \\ \vdots \\ zulu \end{bmatrix}$ 10000

a
aaron
September
visit
zulu

a
aaron
is
visits
zulu

a
⋮
zulu

10,000

$y^{(1)}, y^{(2)}$

reject the september as the candidate of 1st word.

in
$y^{<1>}$    $\hat{y}^{<2>}$    $P(y^{<2>}|x, "in")$

$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow$
$x^{<1>}$    $x^{<T_x>}$    in

$P(y^{<1>}, y^{<2>}|x) = P(y^{<1>}|x)\, P(y^{<2>}|x, y^{<1>})$

jane    $\hat{y}^{<2>}$

$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow$
$x^{<1>}$    $x^{<T_x>}$

$P(y^{<2>}|x, "jane")$

Because the beam width is 3, at every step you instantiate three copies of the network to evaluate these partial sentence fragments and the output.
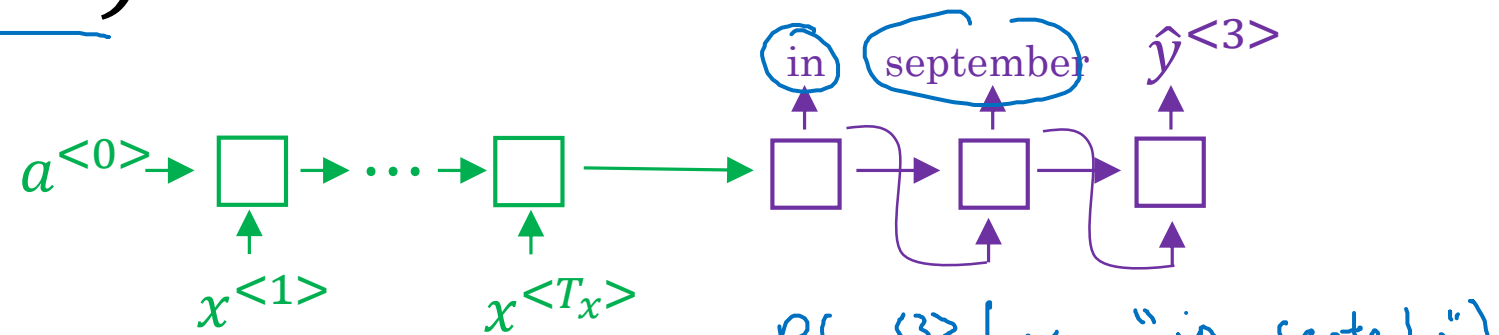
September    $\hat{y}^{<2>}$

$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow$
$x^{<1>}$    $x^{<T_x>}$

You only need three copies of the network to very quickly evaluate all 10,000 possible outputs and that softmax output say y2.

Andrew Ng

# Beam search ($B = 3$)

in september

a
aaron
jane
zulu

jane is

a
visits
zulu

jane visits

a
africa
zulu

$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow \square \rightarrow \square \rightarrow \square$

$x^{<1>} \qquad x^{<T_x>}$

in   september   $\hat{y}^{<3>}$

$P(y^{<3>} \mid x, \text{"in september"})$

jane   is   $\hat{y}^{<3>}$

jane   visits   $\hat{y}^{<3>}$

$P(y^{<1>}, y^{<2>} \mid x)$

jane visits africa in september. <EOS>

Andrew Ng