Examples: Named Entity Recognition

1) Benedict lives in Easthampton, Massachusetts.

Classify each word as a person, location, or other
0                1          2

$x^{<1>}$ = "Benedict", $x^{<2>}$ = "lives", $x^{<3>}$ = "in", $x^{<4>}$ = "Easthampton", $x^{<5>}$ = "Massachusetts"

$y^{<1>}$ = 0    $y^{<2>}$ = 2    $y^{<3>}$ = 2    $y^{<4>}$ = 1    $y^{<5>}$ = 1

$T_x = 5$, $T_y = 5$

2) This movie was the best!

Classify whole sentence as positive or negative
1              0

$T_x = 5$,    $T_y = 1$

$x^{<1>}$ = "This", $\cdots$, $x^{<5>}$ = "best"    $y^{<1>}$ = 1

Order matters!

"The cat chased the mouse."

vs.

"The mouse chased the cat."

3) Translation:

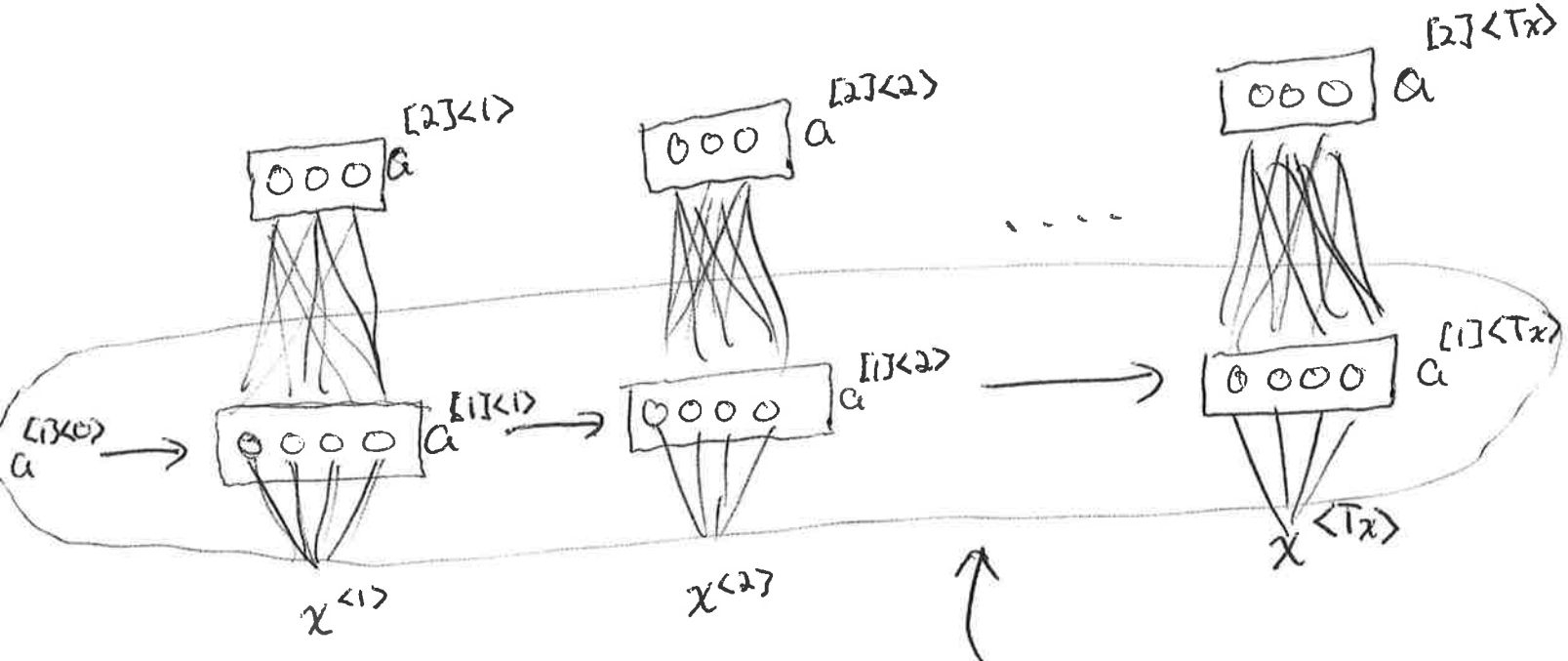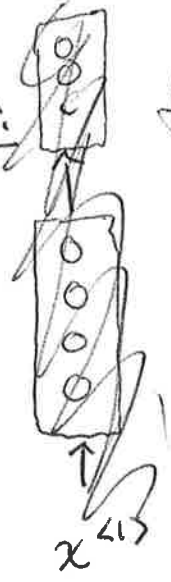Why is the cat so cute? $\rightarrow$ Porqué es el gato tan lindo?
$T_x = 6$                                      $T_y = 7$

# Notation:

$a^{(i)[l]<t>}$: activation in layer $l$ at time $t$ for observation $i$.

## Picture: for named entity recognition example



$x^{<1>}$



$a^{[2]<1>}$  $a^{[2]<2>}$  $a^{[2]<Tx>}$

. . . .

$a^{[3]<0>}$  $a^{[1]<1>}$  $a^{[1]<2>}$  $a^{[1]<Tx>}$

$x^{<1>}$  $x^{<2>}$  $x^{<Tx>}$

a NN that takes in
first word predicts
type of word
(person, location, other)

this (circled) is a
recurrent layer

# Forward Prop.:

$$a^{[1]<t>} = g^{[1]}\left( \left(W_{aa}^{[1]}\right)^T \cdot a^{[1]<t-1>} + \left(W_{ax}^{[1]}\right)^T x^{<t>} + b_a^{[1]} \right)$$

$$= g^{[1]}\left( \left[ \left(W_{aa}^{[1]}\right)^T \left(W_{ax}^{[1]}\right)^T \right] \begin{bmatrix} a^{[1]<t-1>} \\ x^{<t>} \end{bmatrix} + b_a^{[1]} \right)$$

For example suppose we have 4 units in each box of the recurrent layer, and 100 input features.

$W_{aa}$ is $4 \times 4$

$W_{ax}$ is $100 \times 4$

$b_a$ is $4 \times 1$

For recurrent layers, $g$ is almost always tanh activation (ReLU also occasionally used)
  ↳ prevents exploding gradients
  ↳ vanishing gradients is a serious problem we will address through other strategies.

In our example,
$$a^{[2]<t>} = g^{[2]}\left( \left(W_{aa}^{[2]}\right)^T a^{[1]<t>} + b^{[2]} \right)$$

$g^{[2]}$ is whatever appropriate activation for your task
(eg. softmax for named entity recognition with 3 classes)

# Loss function:

$$J(b, \omega) = \frac{1}{T_y} \sum_{i=1}^{m} \sum_{t=1}^{T_y^{(i)}} J^{(i)<t>}(b, \omega)$$

basically, add up losses across all subjects and time points.

Formally, assumes all examples are independent. $(i=1, \ldots, m)$

Within one example, allows for dependence; sort of.

Distribution of $y^{(i)<t>}$ makes use of $x^{(2)1>}, \ldots, x^{(i)<t>}$

Assumes conditional independence:

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

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

$$\cdots \cdot P(y^{<t>} = y^{<t>} \mid x^{<1>}, x^{<2>}, \ldots, x^{<t>})$$

"If I know all of the inputs up to the current time, $x^{<1>}, \ldots, x^{<t>}$, knowing $y^{<t-1>}$ would not give me any additional information about $y^{<t>}$."

↳ not realistic, but better than ignoring time.

# Backward Propagation Through Time.

Start with the cost function <u>at the last time point.</u>

$$\frac{\partial J}{\partial w} = \sum_{i=1}^{m} \sum_{t=1}^{T_y^{(i)}} \frac{\partial}{\partial w} J^{(i)<t>}$$

$$\frac{\partial J^{(i)<t>}}{\partial w} = \frac{\partial J^{(i)<t>}}{\partial a^{[2]<t>}} \cdot \frac{\partial a^{[2]<t>}}{\partial a^{[1]<t>}} \cdot \frac{\partial a^{[1]<t>}}{\partial w}$$

$$a^{[1]<t>} = g^{[1]}\left( (w^{[1]})^T \begin{bmatrix} a^{[1]<t-1>} \\ x^{<t>} \end{bmatrix} + b^{[1]} \right)$$

So

$$\frac{\partial}{\partial w} a^{[1]<t>} = \frac{\partial a^{[1]<t>}}{\partial z^{[1]<t>}} \cdot \frac{\partial z^{[1]<t>}}{\partial w}$$

$$= \frac{\partial a^{[1]<t>}}{\partial z^{[1]<t>}} \left( \frac{\partial}{\partial w} w_a^{[1]T} a^{[1]<t-1>} + \frac{\partial}{\partial w} (w_x^{[1]})^T x^{<t>} \right)$$

w also need to calculate $a^{[1]<t-1>}$ .