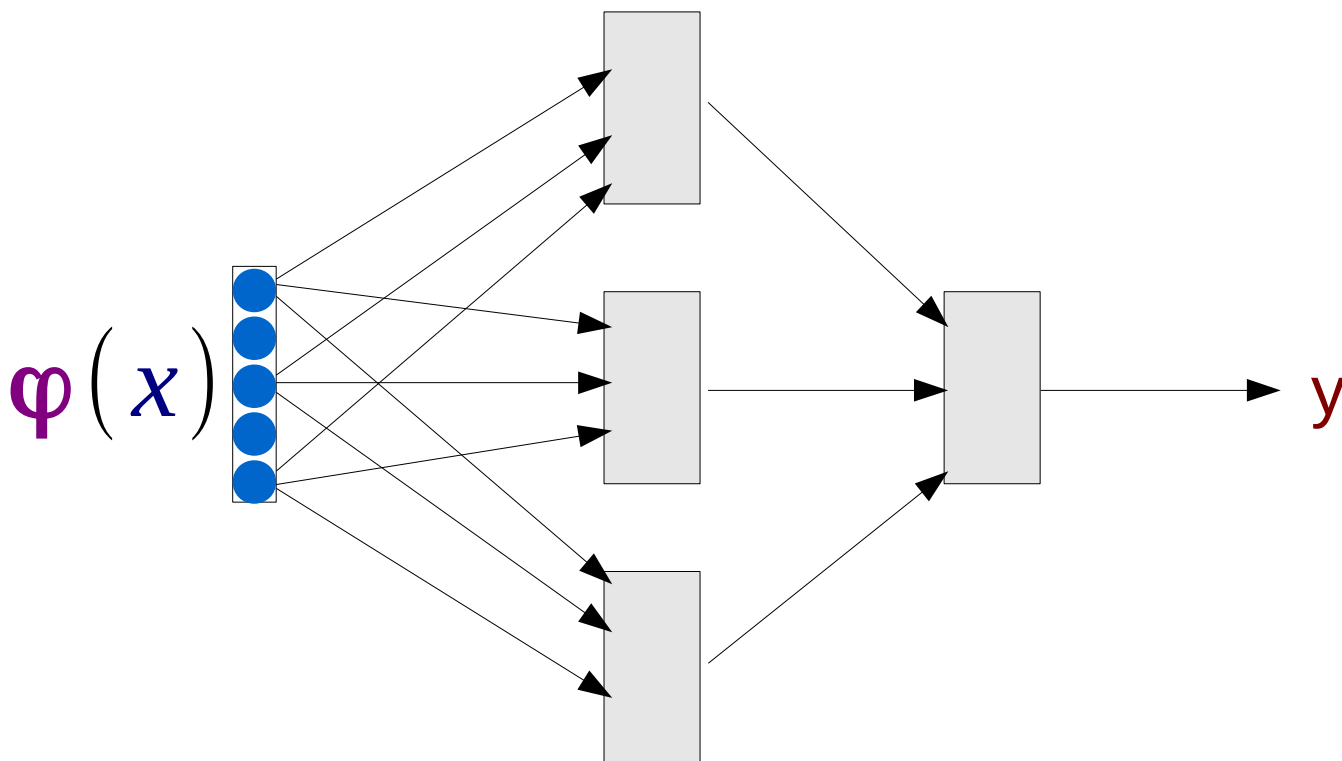


NLP Programming Tutorial 8 - Recurrent Neural Nets

Graham Neubig
Nara Institute of Science and Technology (NAIST)

Feed Forward Neural Nets

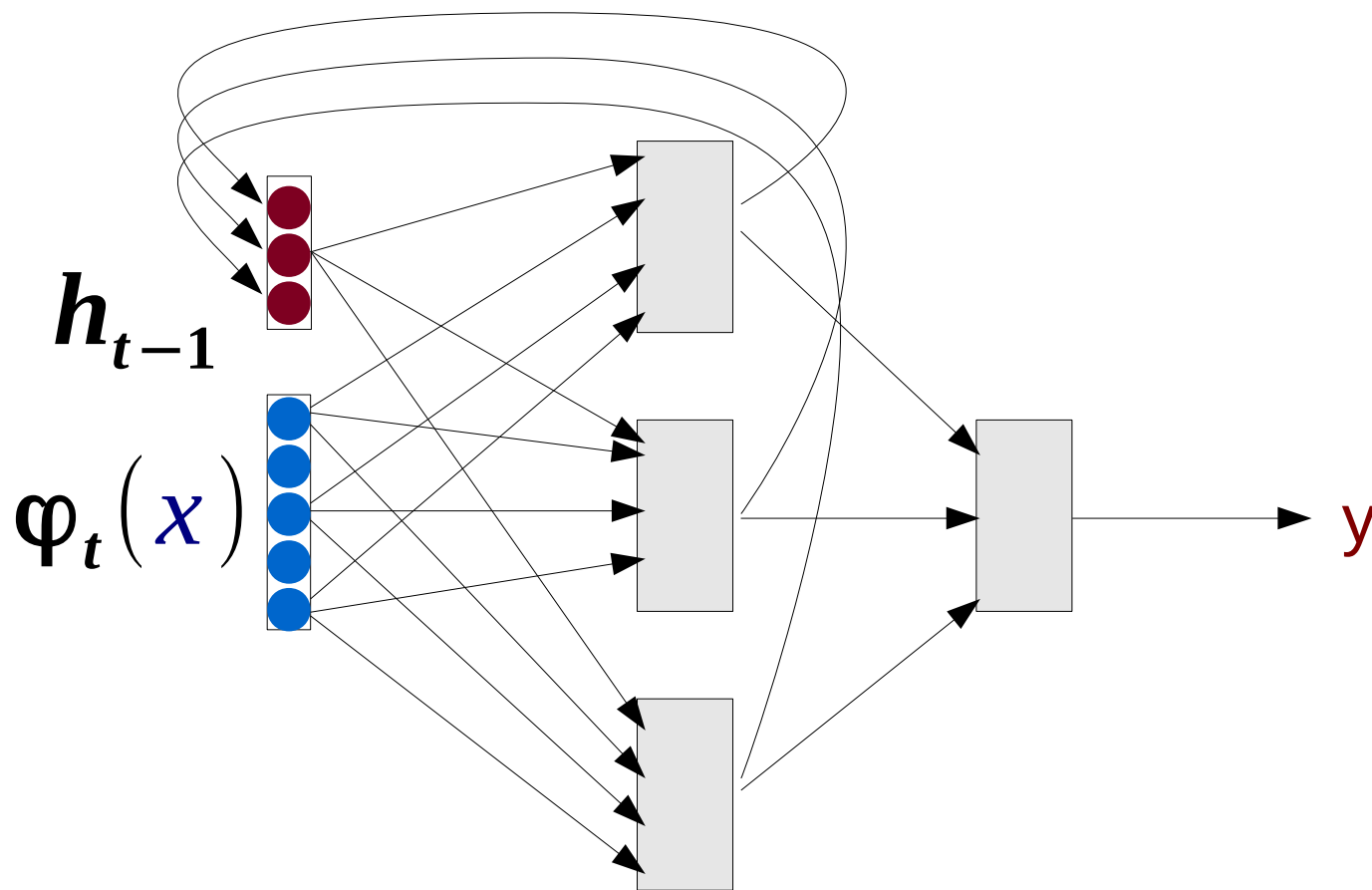
- All connections point forward



- It is a directed acyclic graph (DAG)

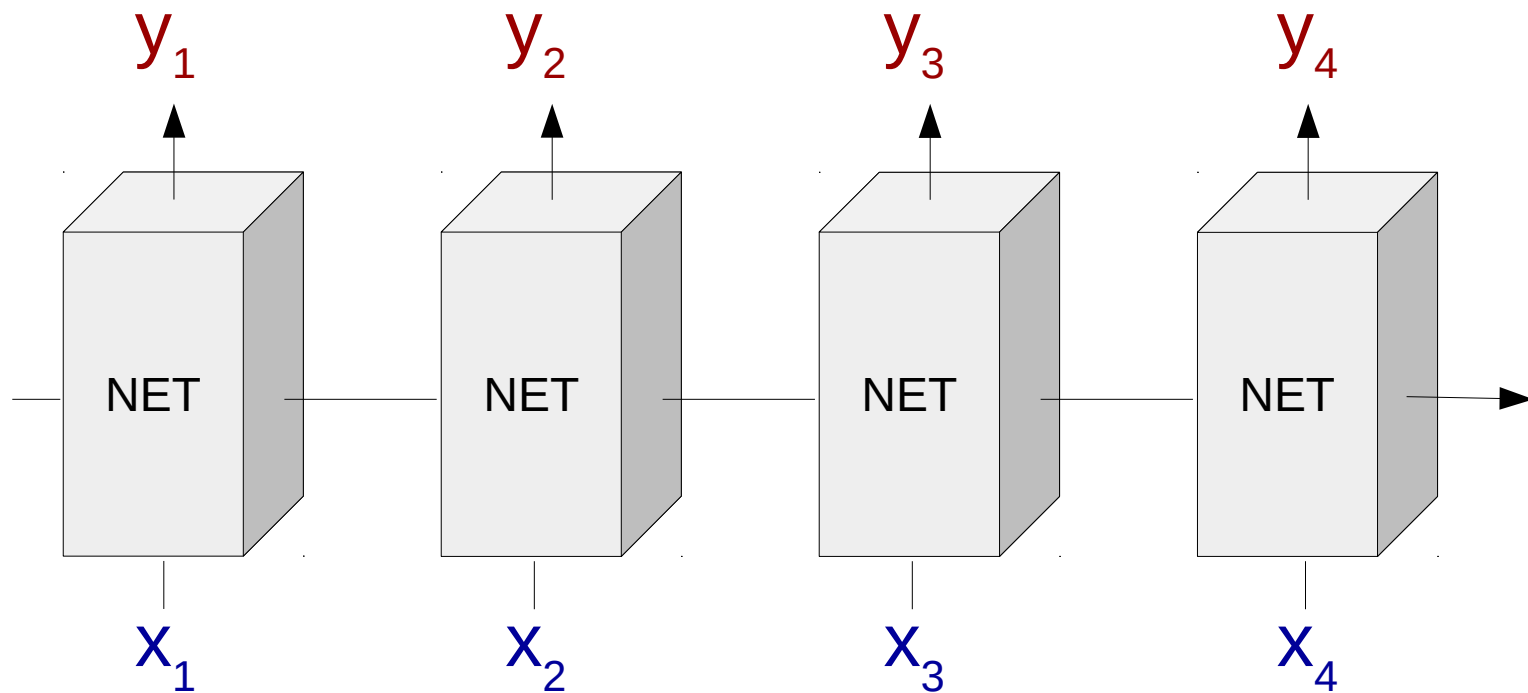
Recurrent Neural Nets (RNN)

- Part of the node outputs return as input

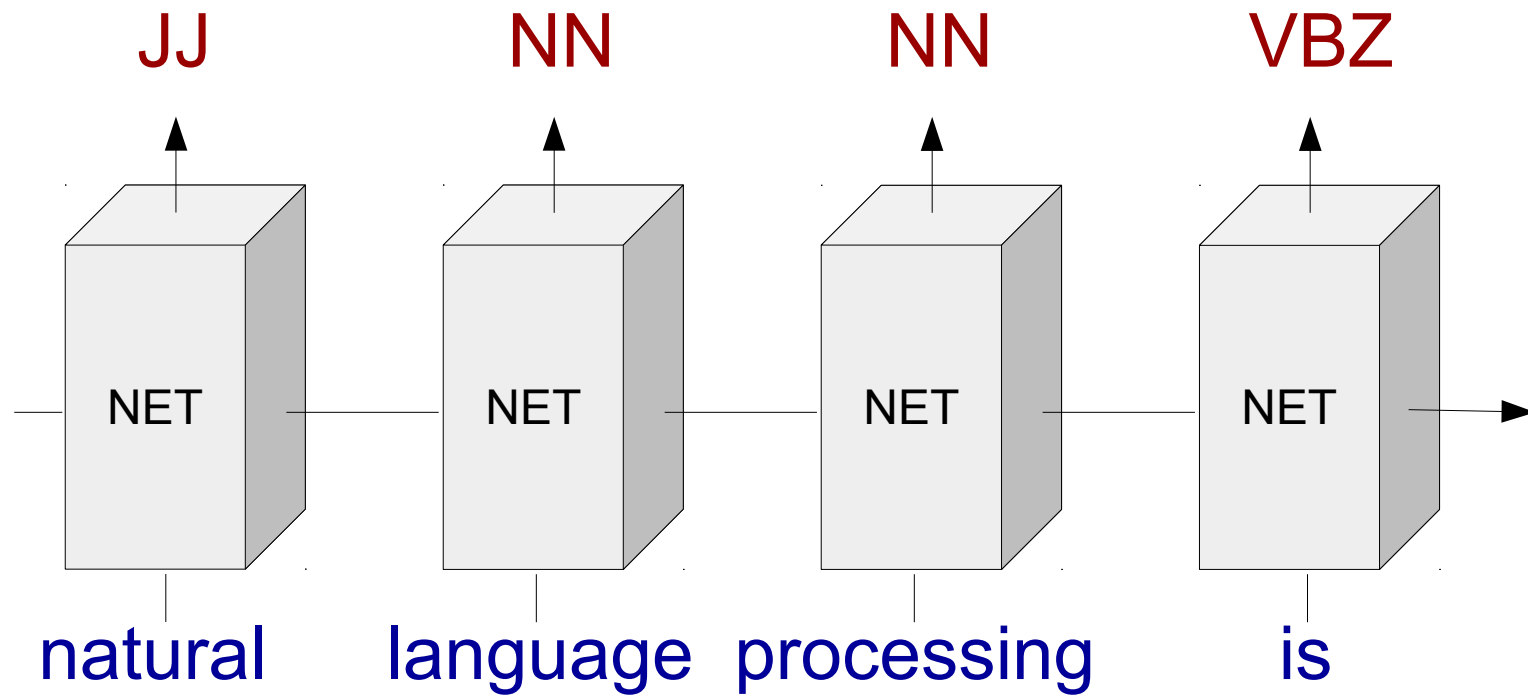


- Why? It is possible to “memorize”

RNN in Sequence Modeling



Example: POS Tagging



Multi-class Prediction with Neural Networks

Review: Prediction Problems

Given x ,

predict y

A book review

Oh, man I love this book!
This book is so boring...

Is it positive?

yes
no

Binary
Prediction
(2 choices)

A tweet

On the way to the park!
公園に行くなう！

Its language

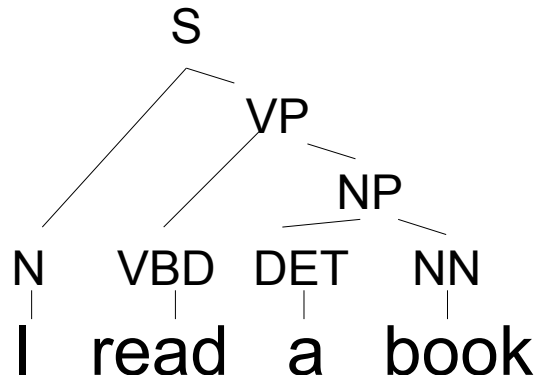
English
Japanese

Multi-class
Prediction
(several choices)

A sentence

I read a book

Its syntactic parse



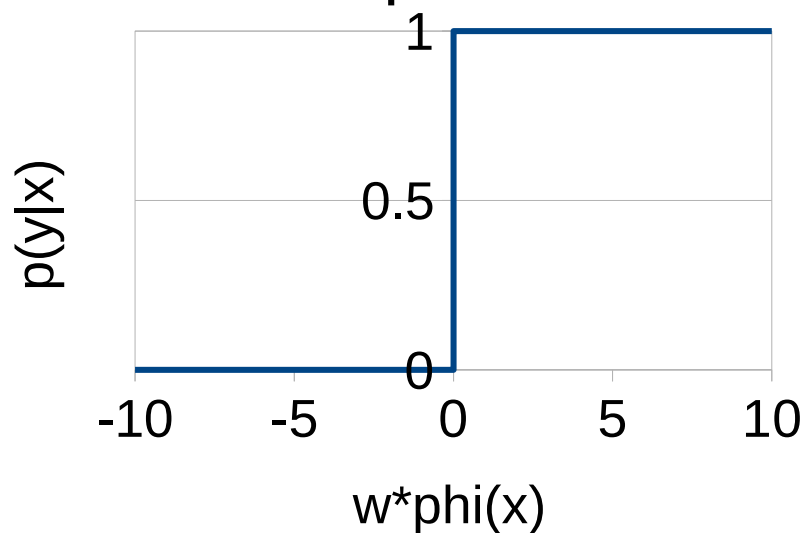
Structured
Prediction
(millions of choices)

Review: Sigmoid Function

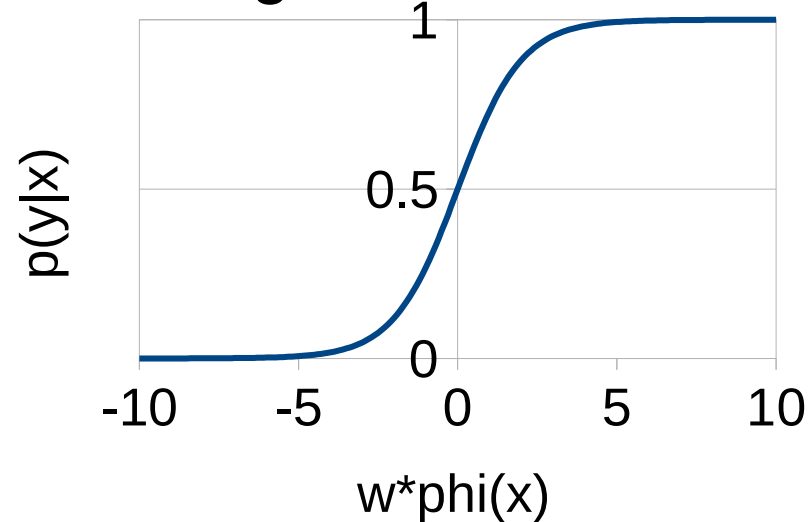
- The sigmoid softens the step function

$$P(y = 1 | x) = \frac{e^{w \cdot \varphi(x)}}{1 + e^{w \cdot \varphi(x)}}$$

Step Function



Sigmoid Function



softmax Function

- Sigmoid function for multiple classes

$$P(y|x) = \frac{e^{w \cdot \varphi(x, y)}}{\sum_{\tilde{y}} e^{w \cdot \varphi(x, \tilde{y})}}$$

← Current class

← Sum of other classes

- Can be expressed using matrix/vector ops

$$\mathbf{r} = \exp(\mathbf{W} \cdot \boldsymbol{\varphi}(\mathbf{x}))$$

$$\mathbf{p} = \mathbf{r} / \sum_{\tilde{r} \in r} \tilde{r}$$

Selecting the Best Value from a Probability Distribution

- Find the index y with the highest probability

```
find_best( $p$ ):  
   $y = 0$   
  for each element  $i$  in  $1 .. \text{len}(p)-1$ :  
    if  $p[i] > p[y]$ :  
       $y = i$   
  return  $y$ 
```

softmax Function Gradient

- The difference between the true and estimated probability distributions

$$-d \text{err} / d \varphi_{out} = \mathbf{p}' - \mathbf{p}$$

- The true distribution \mathbf{p}' is expressed with a vector with only the y -th element 1 (a one-hot vector)

$$\mathbf{p}' = \{0, 0, \dots, 1, \dots, 0\}$$

Creating a 1-hot Vector

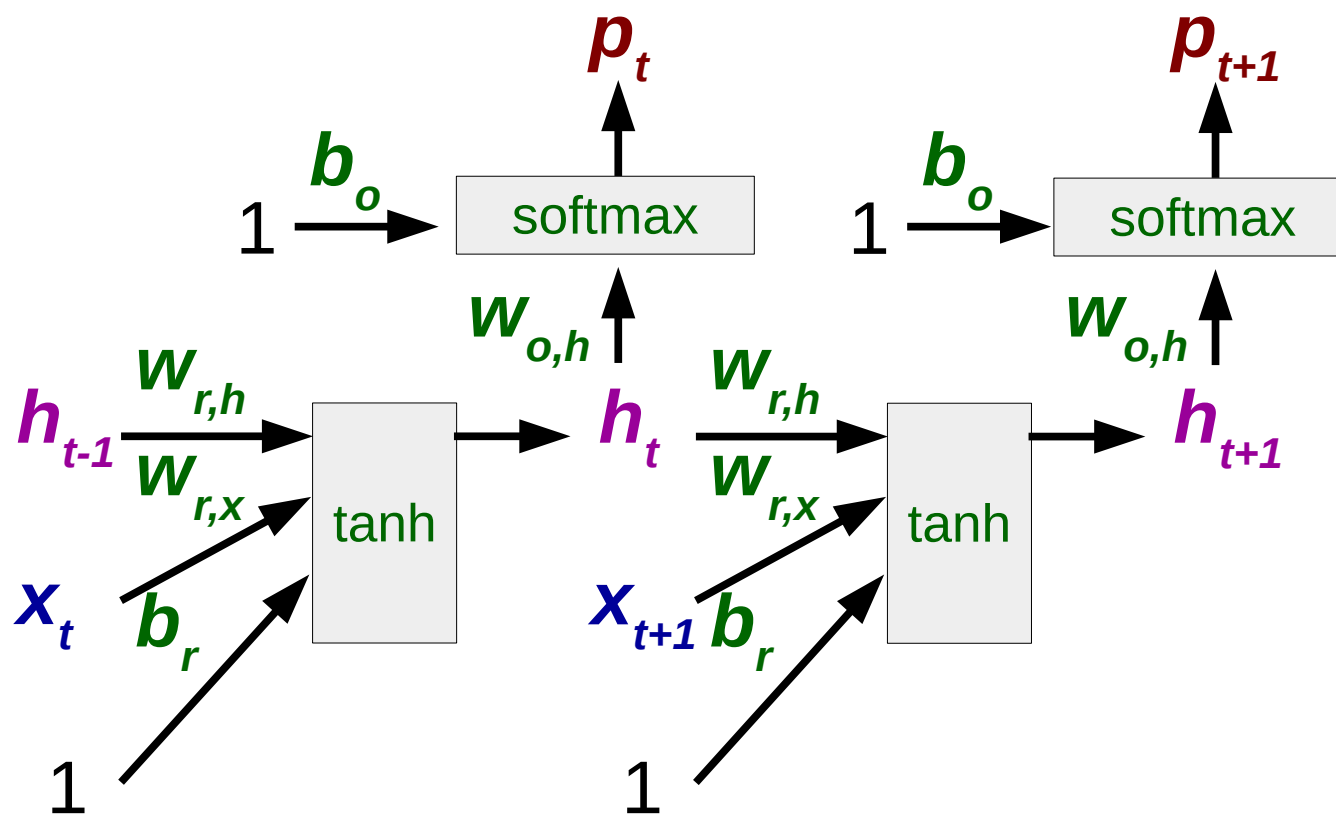
```
create_one_hot(id, size):  
    vec = np.zeros(size)  
    vec[id] = 1  
    return vec
```

Forward Propagation in Recurrent Nets

Review: Forward Propagation Code

```
forward_nn(network,  $\varphi_0$ )  
   $\varphi$  = [  $\varphi_0$  ] # Output of each layer  
  for each layer i in 1 .. len(network):  
     $w$ ,  $b$  = network[i-1]  
    # Calculate the value based on previous layer  
     $\varphi[i]$  = np.tanh( np.dot(  $w$ ,  $\varphi[i-1]$  ) +  $b$  )  
  return  $\varphi$  # Return the values of all layers
```

RNN Calculation



$$h_t = \tanh(w_{r,h} \cdot h_{t-1} + w_{r,x} \cdot x_t + b_r)$$

$$p_t = \text{softmax}(w_{o,h} \cdot h_t + b_o)$$

RNN Forward Calculation

```
forward_rnn( $w_{r,x}$ ,  $w_{r,h}$ ,  $b_r$ ,  $w_{o,h}$ ,  $b_o$ ,  $x$ )
   $h = []$  # Hidden layers (at time t)
   $p = []$  # Output probability distributions (at time t)
   $y = []$  # Output values (at time t)
  for each time  $t$  in  $0 .. \text{len}(x)-1$ :
    if  $t > 0$ :
       $h[t] = \tanh(w_{r,x}x[t] + w_{r,h}h[t-1] + b_r)$ 
    else:
       $h[t] = \tanh(w_{r,x}x[t] + b_r)$ 
       $p[t] = \tanh(w_{o,h}h[t] + b_o)$ 
       $y[t] = \text{find\_max}(p[t])$ 
  return  $h, p, y$ 
```


Review: Back Propagation in Feed-forward Nets

Stochastic Gradient Descent

- Online training algorithm for probabilistic models (including logistic regression)

$w = 0$

for / iterations

for each labeled pair x , y in the data

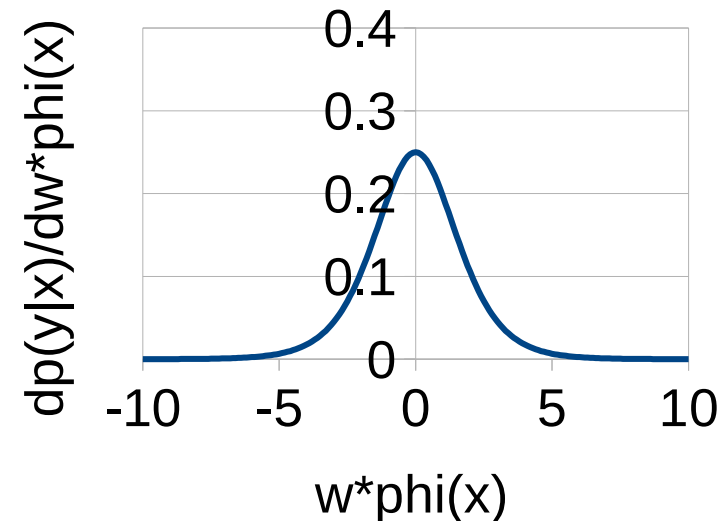
$w += \alpha * dP(y|x)/dw$

- **In other words**
 - For every training example, calculate the **gradient** (the direction that will increase the probability of y)
 - **Move** in that direction, multiplied by learning rate α

Gradient of the Sigmoid Function

- Take the derivative of the probability

$$\begin{aligned} \frac{d}{d w} P(y=1|x) &= \frac{d}{d w} \frac{e^{w \cdot \varphi(x)}}{1+e^{w \cdot \varphi(x)}} \\ &= \varphi(x) \frac{e^{w \cdot \varphi(x)}}{(1+e^{w \cdot \varphi(x)})^2} \end{aligned}$$

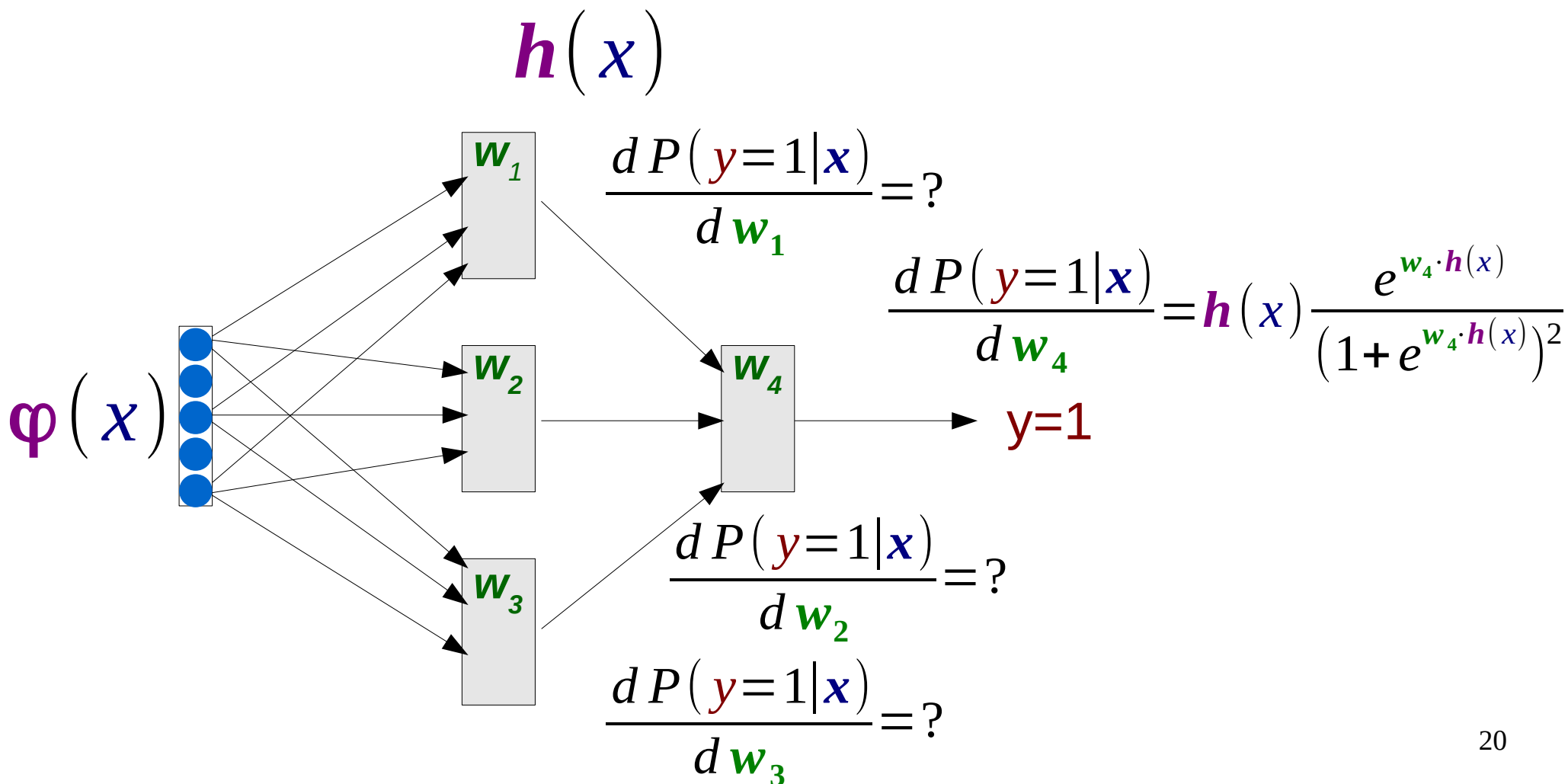


$$\begin{aligned} \frac{d}{d w} P(y=-1|x) &= \frac{d}{d w} \left(1 - \frac{e^{w \cdot \varphi(x)}}{1+e^{w \cdot \varphi(x)}} \right) \\ &= -\varphi(x) \frac{e^{w \cdot \varphi(x)}}{(1+e^{w \cdot \varphi(x)})^2} \end{aligned}$$

Learning:

Don't Know Derivative for Hidden Units!

- For NNs, only know correct tag for last layer



Answer: Back-Propogation

- Calculate derivative w/ chain rule

$$\frac{dP(y=1|x)}{dw_1} = \frac{dP(y=1|x)}{dw_4 h(x)} \frac{dw_4 h(x)}{dh_1(x)} \frac{dh_1(x)}{dw_1}$$

$$\frac{e^{w_4 \cdot h(x)}}{(1 + e^{w_4 \cdot h(x)})^2}$$

↓

Error of next unit (δ_4)

$$w_{1,4}$$

↓

Weight

$$\frac{dh_1(x)}{dw_1}$$

↓

Gradient of this unit

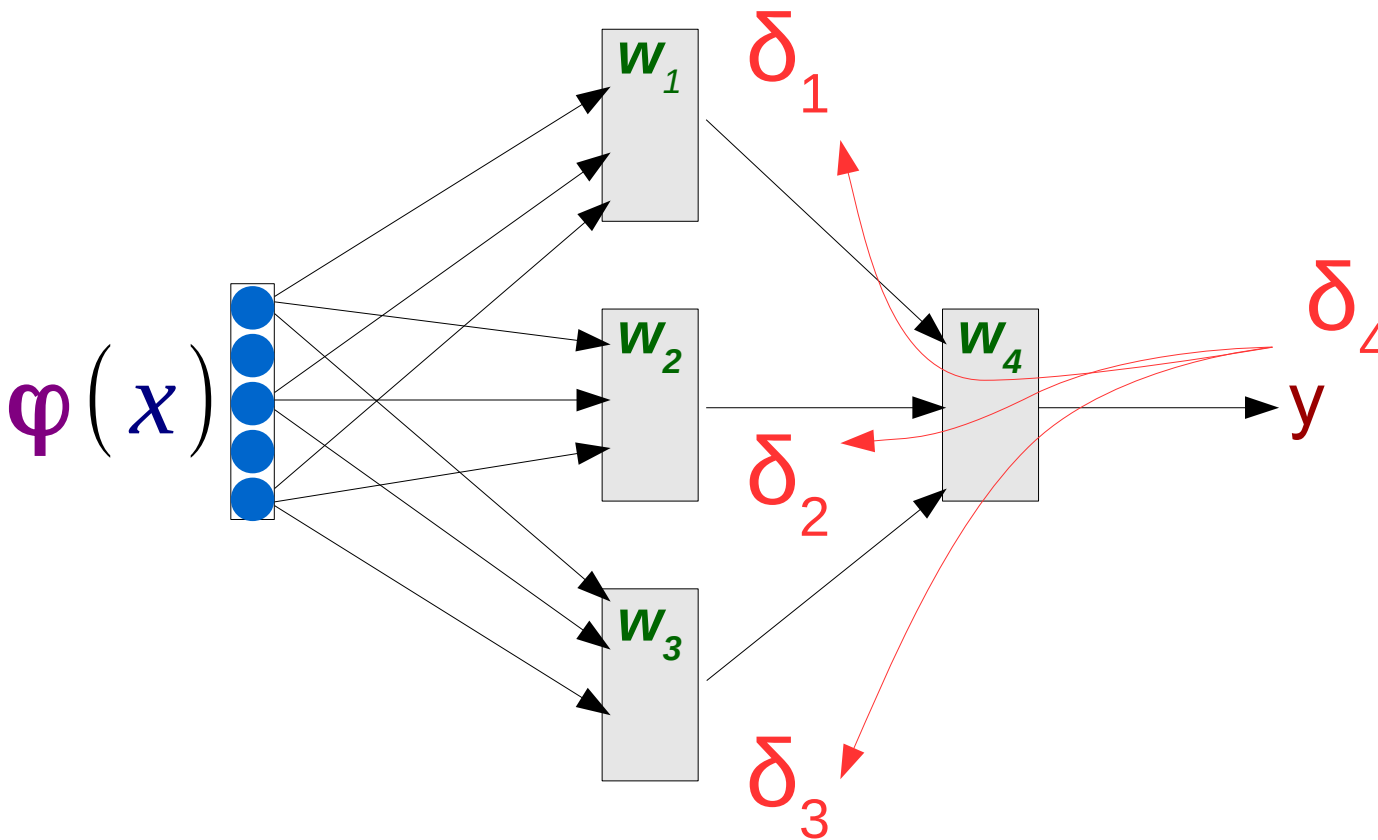
In General

Calculate i based on next units j :

$$\frac{dP(y=1|x)}{dw_i} = \frac{dh_i(x)}{dw_i} \sum_j \delta_j w_{i,j}$$

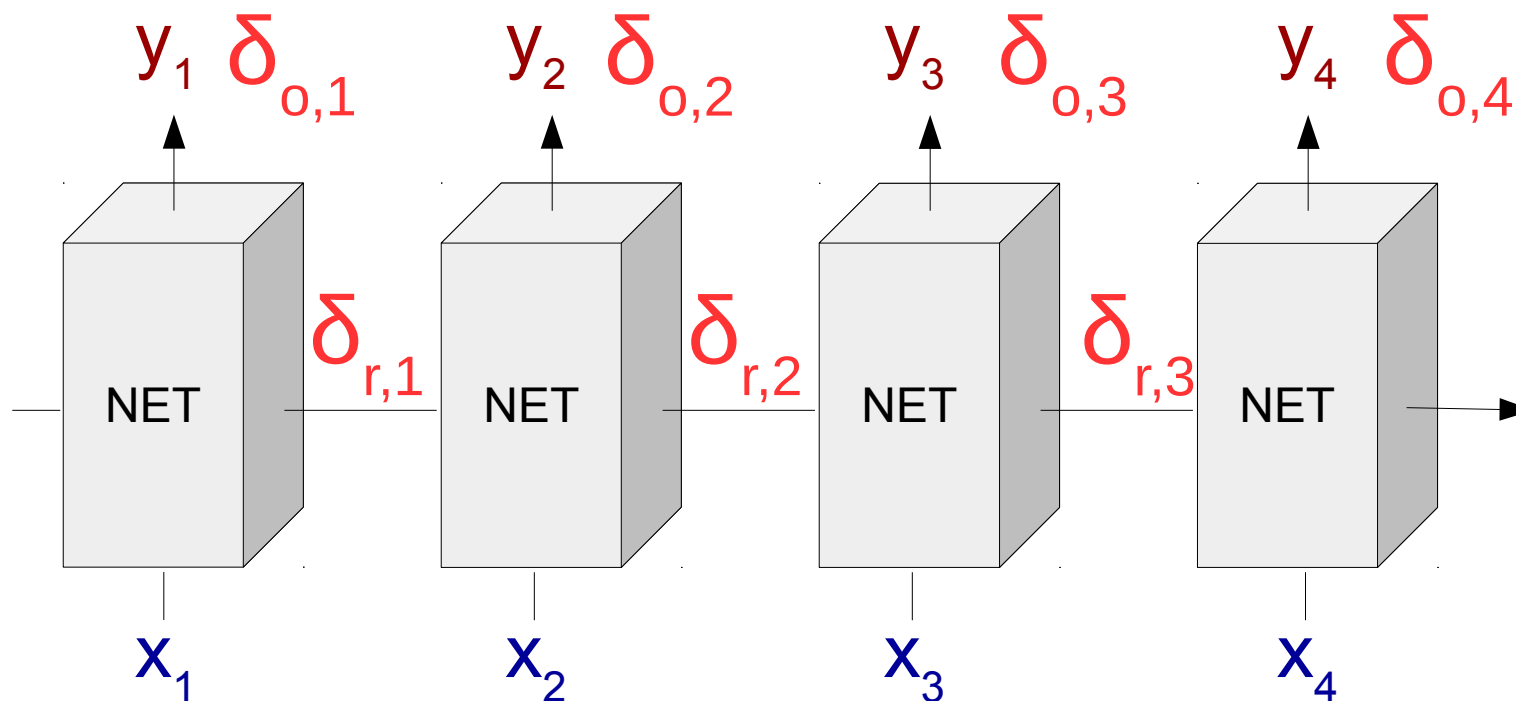
Conceptual Picture

- Send errors back through the net



Back Propagation in Recurrent Nets

What Errors do we Know?

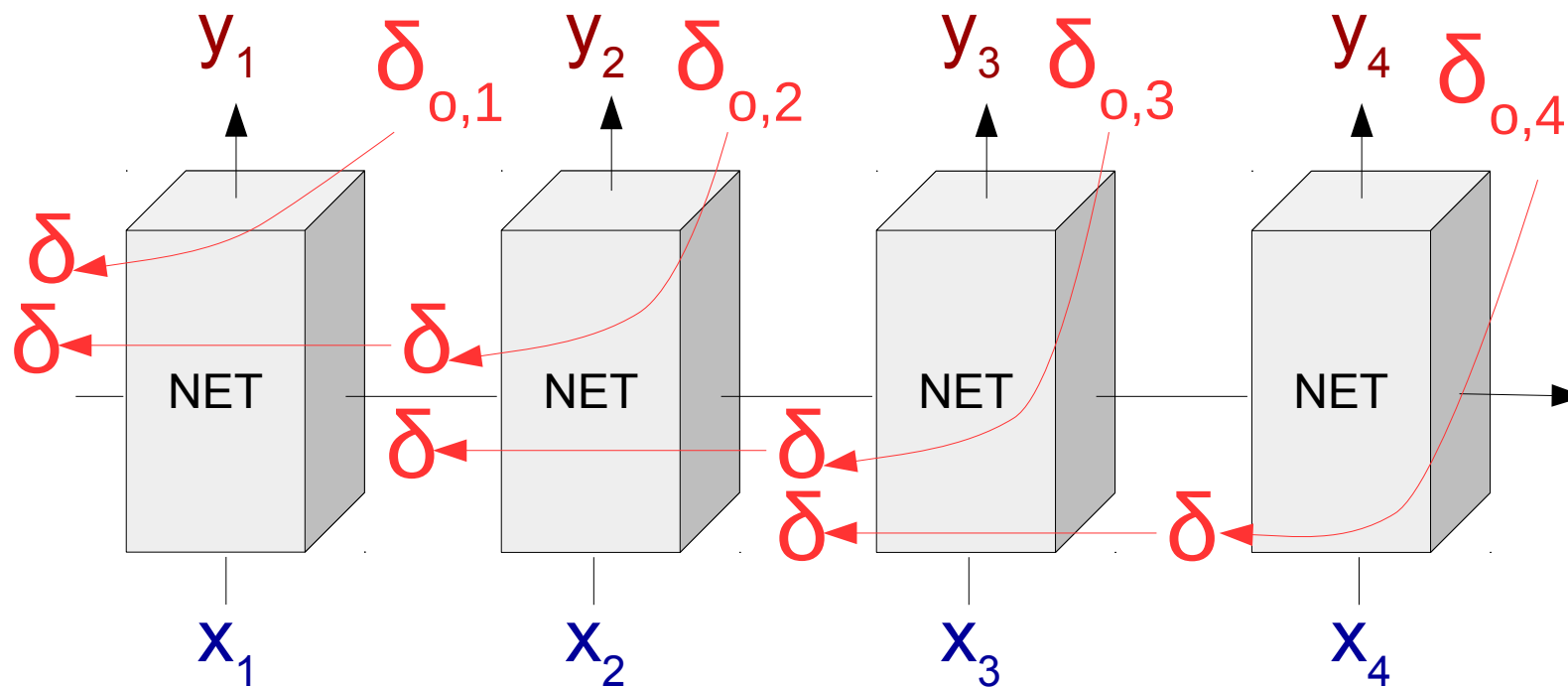


- We know the output errors δ_o
- Must use back-prop to find recurrent errors δ_r

How to Back-Propagate?

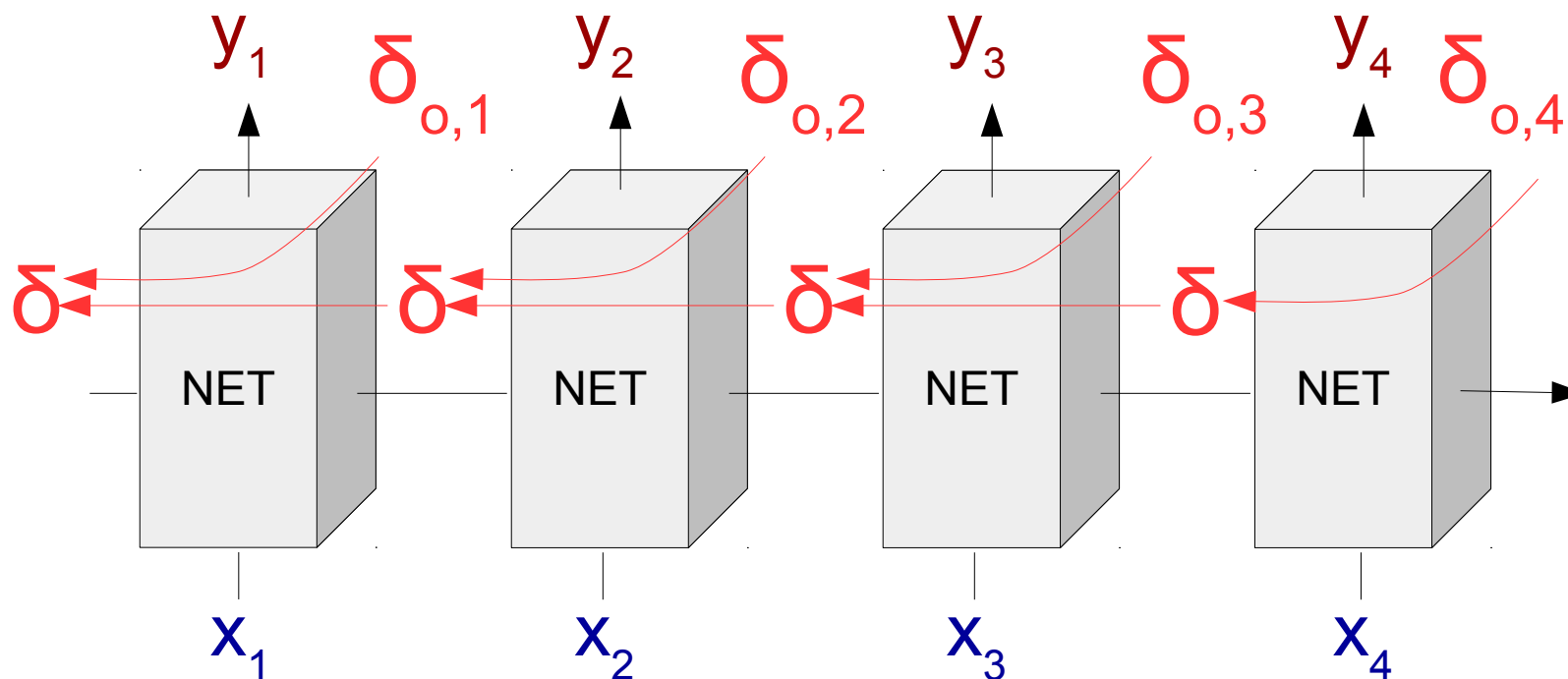
- Standard **back propagation through time (BPTT)**
 - For each δ_o , calculate n steps of δ_r
- **Full gradient** calculation
 - Use dynamic programming to calculate the whole sequence

Back Propagation through Time



- Use only one output error
- Stop after n steps (here, n=2)

Full Gradient Calculation

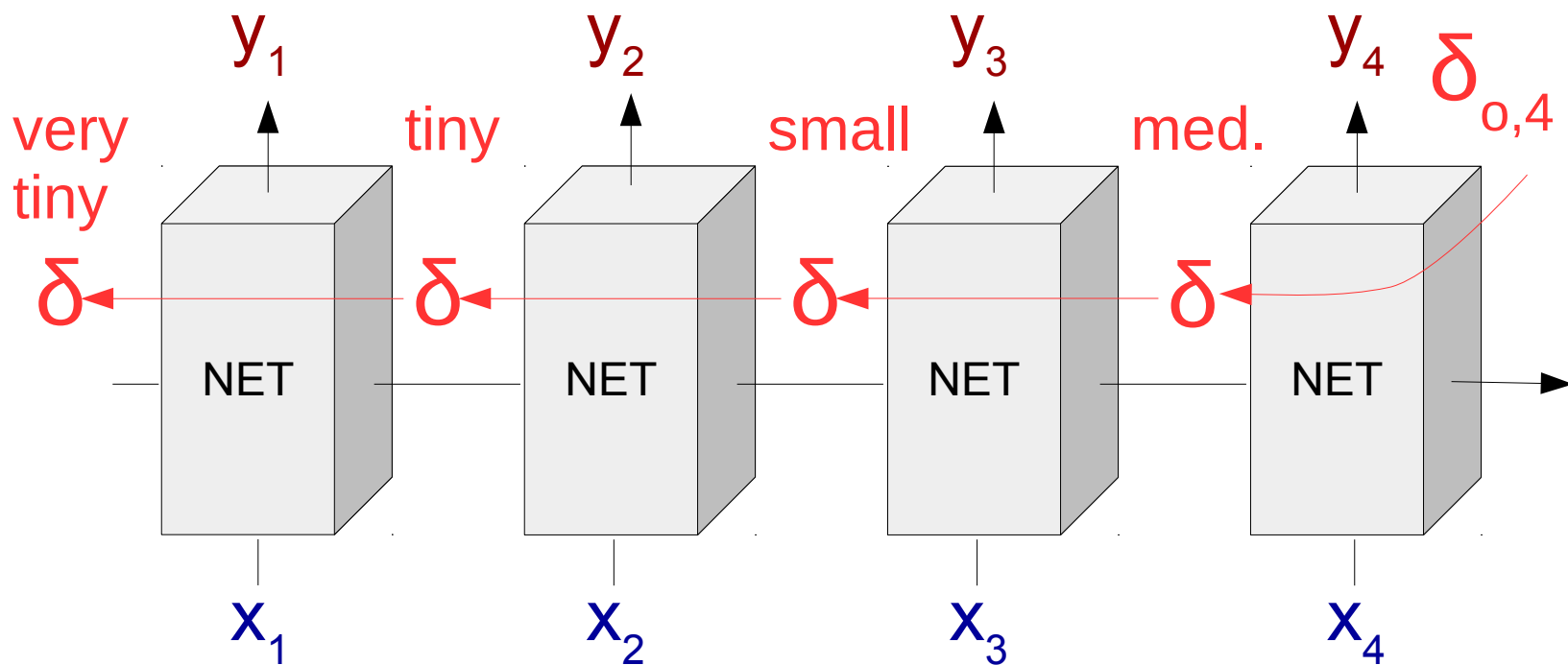


- First, calculate whole net result forward
- Then, calculate result backwards

BPTT? Full Gradient?

- Full gradient:
 - + Faster, no time limit
 - - Must save the result of the whole sequence in memory
- BPTT:
 - + Only remember the results in the past few steps
 - - Slower, less accurate for long dependencies

Vanishing Gradient in Neural Nets



- “Long Short Term Memory” is designed to solve this

RNN Full Gradient Calculation

```

gradient_rnn( $w_{r,x}$ ,  $w_{r,h}$ ,  $b_r$ ,  $w_{o,h}$ ,  $b_o$ ,  $x$ ,  $h$ ,  $p$ ,  $y'$ )
initialize  $\Delta w_{r,x}$ ,  $\Delta w_{r,h}$ ,  $\Delta b_r$ ,  $\Delta w_{o,h}$ ,  $\Delta b_o$ 
 $\delta_r'$  = np.zeros(len( $b_r$ )) # Error from the following time step
for each time  $t$  in len( $x$ )-1 .. 0:
     $p'$  = create_one_hot( $y'[t]$ )
     $\delta_o'$  =  $p' - p[t]$  # Output error
     $\Delta w_{o,h}$  += np.outer( $h[t]$ ,  $\delta_o'$ );     $\Delta b_o$  +=  $\delta_o'$  # Output gradient
     $\delta_r$  = np.dot( $\delta_r'$ ,  $w_{r,h}$ ) + np.dot( $\delta_o'$ ,  $w_{o,h}$ ) # Backprop
     $\delta_r' = \delta_r * (1 - h[t]^2)$  # tanh gradient
     $\Delta w_{r,x}$  += np.outer( $x[t]$ ,  $\delta_r'$ );     $\Delta b_r$  +=  $\delta_r'$  # Hidden gradient
    if  $t \neq 0$ :
         $\Delta w_{r,h}$  += np.outer( $h[t-1]$ ,  $\delta_r'$ );
return  $\Delta w_{r,x}$ ,  $\Delta w_{r,h}$ ,  $\Delta b_r$ ,  $\Delta w_{o,h}$ ,  $\Delta b_o$ 
    
```

Weight Update

```
update_weights( $w_{r,x}$ ,  $w_{r,h}$ ,  $b_r$ ,  $w_{o,h}$ ,  $b_o$ ,  $\Delta w_{r,x}$ ,  $\Delta w_{r,h}$ ,  $\Delta b_r$ ,  $\Delta w_{o,h}$ ,  $\Delta b_o$ ,  $\lambda$ )  
 $w_{r,x} += \lambda * \Delta w_{r,x}$   
 $w_{r,h} += \lambda * \Delta w_{r,h}$   
 $b_r += \lambda * \Delta b_r$   
 $w_{o,h} += \lambda * \Delta w_{o,h}$   
 $b_o += \lambda * \Delta b_o$ 
```

Overall Training Algorithm

```
# Create features
create map x_ids, y_ids, array data
for each labeled pair x, y in the data
    add (create_ids(x, x_ids), create_ids(y, y_ids)) to data
initialize net randomly

# Perform training
for / iterations
    for each labeled pair x, y' in the feat_lab
        h, p, y = forward_rnn(net,  $\varphi_0$ )
         $\Delta$  = gradient_rnn(net, x, h, y')
        update_weights(net,  $\Delta$ ,  $\lambda$ )

print net to weight_file
print x_ids, y_ids to id_file
```


Exercise

Exercise

- Create an RNN for sequence labeling
- Training train-rnn and testing test-rnn
- **Test:** Same data as POS tagging
 - Input: test/05-`{train,test}`-input.txt
 - Reference: test/05-`{train,test}`-answer.txt
- **Train** a model with data/wiki-en-train.norm_pos and **predict** for data/wiki-en-test.norm
- **Evaluate** the POS performance, and **compare** with HMM:
script/gradeupos.pl data/wiki-en-test.pos my_answer.pos