# Vocabulary for Deep Learning (In 40 Minutes! With Pictures!)

Graham Neubig

Nara Institute of Science and Technology (NAIST)

# Prediction Problems

# Given x, predict y
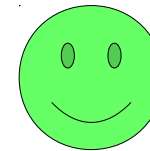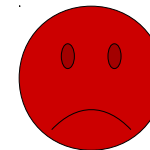
# Example: Sentiment Analysis
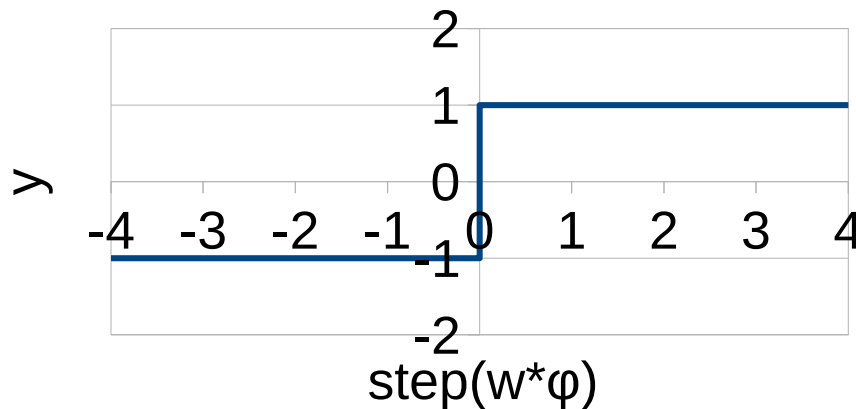
X          Y

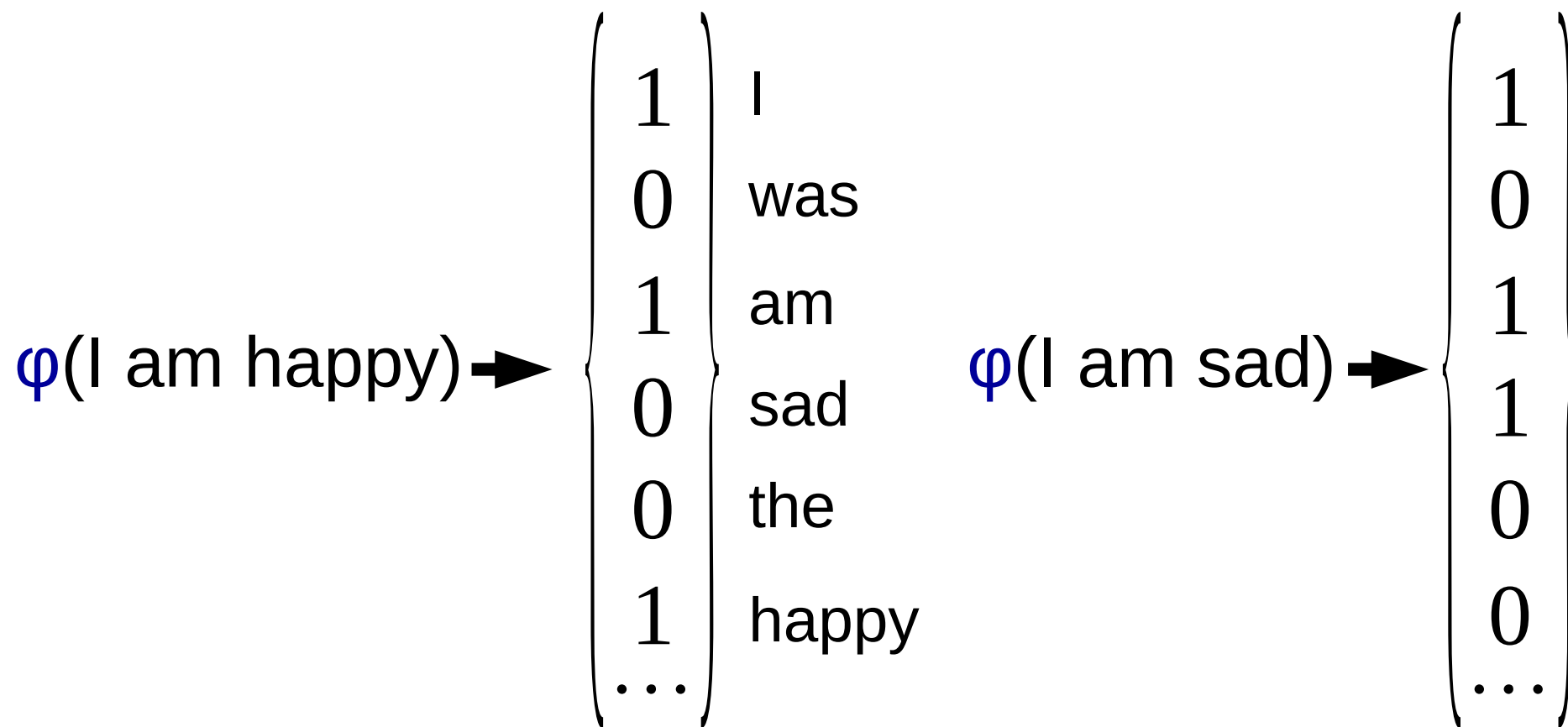I am happy          +1 🙂

I am sad          -1 ☹️

# Linear Classifiers

$$y = \text{step}\left(\boldsymbol{w}^T \boldsymbol{\varphi}(x)\right)$$

- x: the input

- **φ(x)**: vector of feature functions

- w: the weight vector

- y: the prediction, +1 if "yes", -1 if "no"



step(w*φ)

# Example Feature Functions: Unigram Features

$$\varphi(\text{I am happy}) \rightarrow \begin{pmatrix} 1 & \text{I} \\ 0 & \text{was} \\ 1 & \text{am} \\ 0 & \text{sad} \\ 0 & \text{the} \\ 1 & \text{happy} \\ \dots & \end{pmatrix} \qquad \varphi(\text{I am sad}) \rightarrow \begin{pmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 0 \\ 0 \\ \dots \end{pmatrix}$$

# The Perceptron

- Think of it as a "machine" to calculate a weighted sum and insert it into an activation function

$$\varphi(x) \rightarrow \boxed{\begin{array}{l} w \\ \text{step}\left(w^T \varphi(x)\right) \end{array}} \rightarrow y$$

# Sigmoid Function (Logistic Function)

- The sigmoid function is a "softened" version of the step function

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



Step Function

Logistic Function

- Can account for uncertainty

7

- Differentiable

# Logistic Regression

- Train based on conditional likelihood

- Find the parameters **w** that maximize the conditional likelihood of all answers $y_i$ given the example $x_i$

$$\hat{w} = \underset{w}{\mathrm{argmax}} \prod_i P(y_i | x_i ; w)$$

- How do we solve this?

# Stochastic Gradient Descent

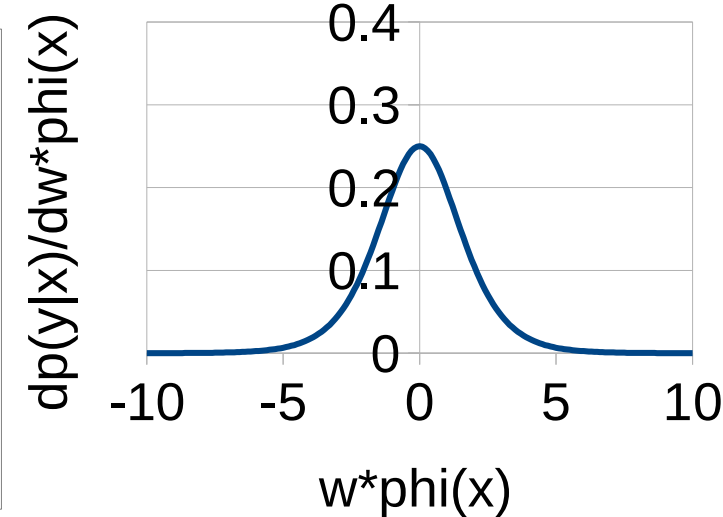- Online training algorithm for probabilistic models (including logistic regression)

**create** map *w*
**for** *I* iterations
    **for each** labeled pair *x, y* in the data
        w += α * 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

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

$$=\; \varphi(x)\frac{e^{w\cdot\varphi(x)}}{(1+e^{w\cdot\varphi(x)})^2}$$

$$\frac{d}{d\,w}P(y=-1|x) \;=\; \frac{d}{d\,w}\left(1-\frac{e^{w\cdot\varphi(x)}}{1+e^{w\cdot\varphi(x)}}\right)$$

$$=\; -\varphi(x)\frac{e^{w\cdot\varphi(x)}}{(1+e^{w\cdot\varphi(x)})^2}$$



10

# Neural Networks

# Problem: Linear Constraint

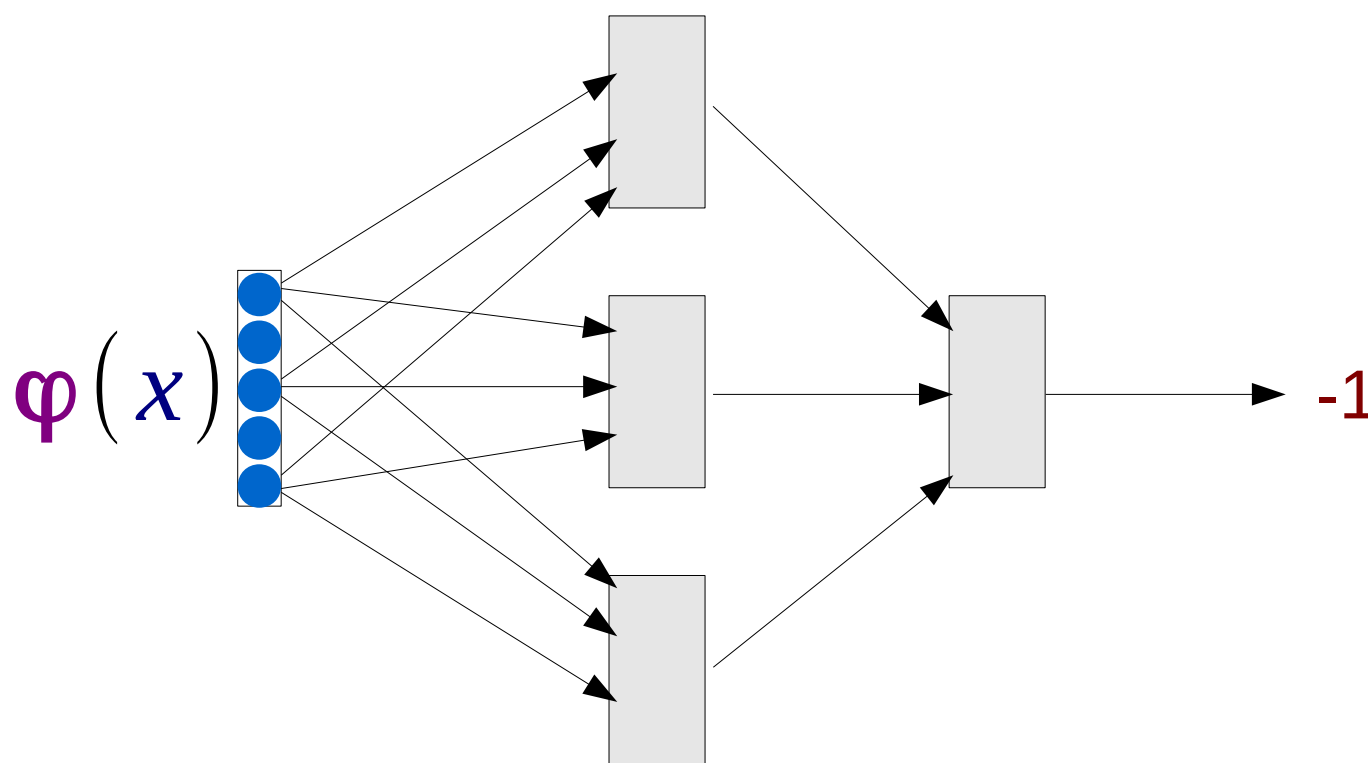- Perceptron cannot achieve high accuracy on non-linear functions

<div align="center">

X      O

O      X

</div>

- Example:  "I am **not** happy"
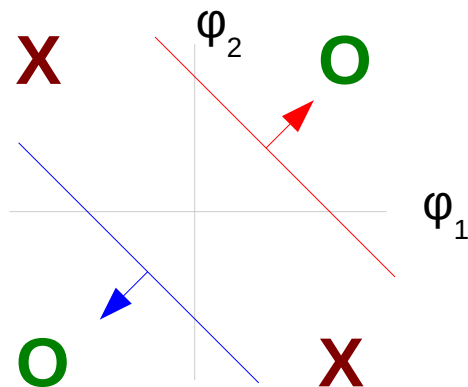
# Neural Networks
# (Multi-Layer Perceptron)

- Neural networks connect multiple perceptrons together
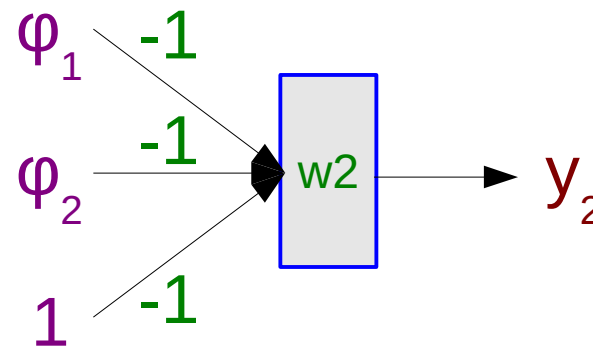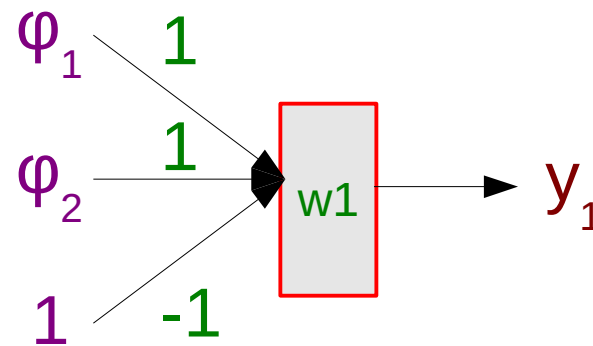
# Example:

- Build two classifiers:

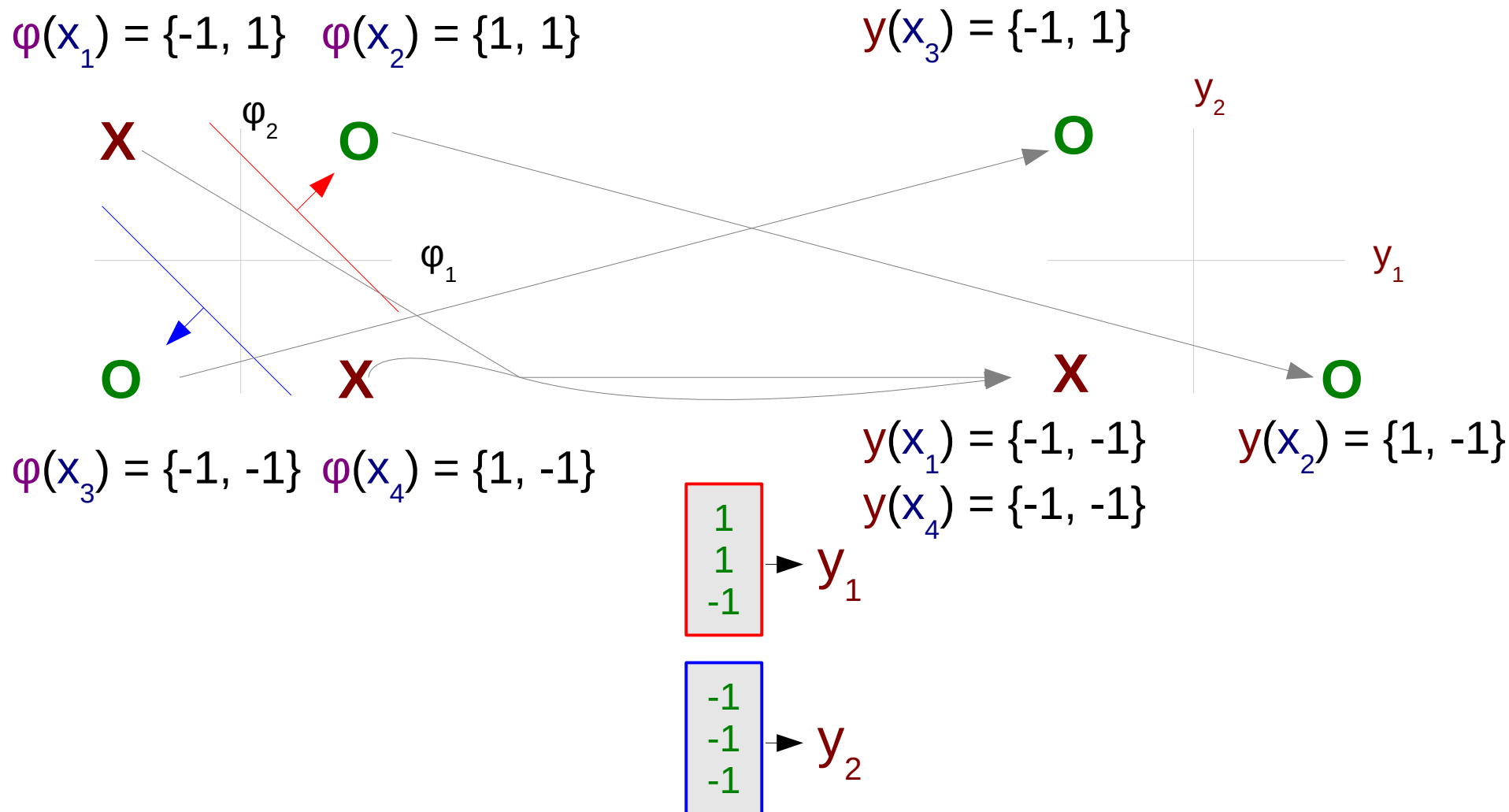$\varphi(x_1) = \{-1, 1\}$   $\varphi(x_2) = \{1, 1\}$
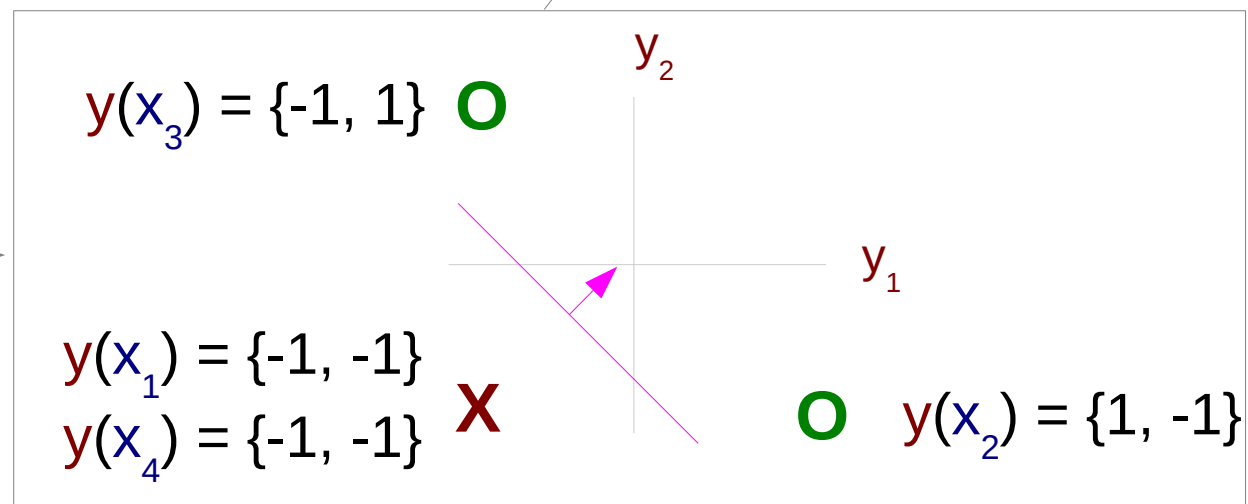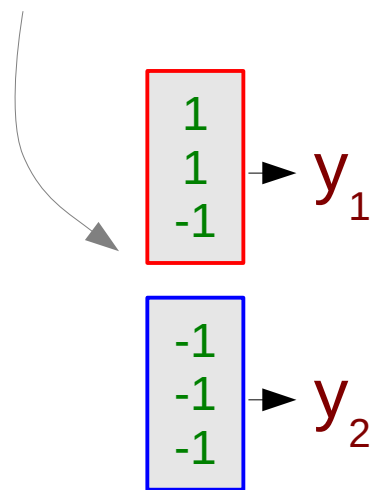


$\varphi(x_3) = \{-1, -1\}$   $\varphi(x_4) = \{1, -1\}$

# Example:

- These classifiers map the points to a new space

$\varphi(x_1) = \{-1, 1\}$   $\varphi(x_2) = \{1, 1\}$         $y(x_3) = \{-1, 1\}$



$\varphi_2$

$\varphi_1$

$y_2$

$y_1$

X    O

O    X

O    X    O

$\varphi(x_3) = \{-1, -1\}$   $\varphi(x_4) = \{1, -1\}$

$y(x_1) = \{-1, -1\}$   $y(x_2) = \{1, -1\}$

$y(x_4) = \{-1, -1\}$

$$\begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix} \rightarrow y_1$$

$$\begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix} \rightarrow y_2$$

15

# Example:

- In the new space, examples are classifiable!

$\varphi(x_1) = \{-1, 1\}$   $\varphi(x_2) = \{1, 1\}$

X   $\varphi_2$   O

$\varphi_1$

O   X

$\varphi(x_3) = \{-1, -1\}$   $\varphi(x_4) = \{1, -1\}$

$\begin{array}{c} 1 \\ 1 \\ 1 \end{array} \rightarrow y_3$

$\begin{array}{c} 1 \\ 1 \\ -1 \end{array} \rightarrow y_1$

$\begin{array}{c} -1 \\ -1 \\ -1 \end{array} \rightarrow y_2$

$y(x_3) = \{-1, 1\}$   O

$y_2$

$y_1$

$y(x_1) = \{-1, -1\}$
$y(x_4) = \{-1, -1\}$   X

O   $y(x_2) = \{1, -1\}$

16
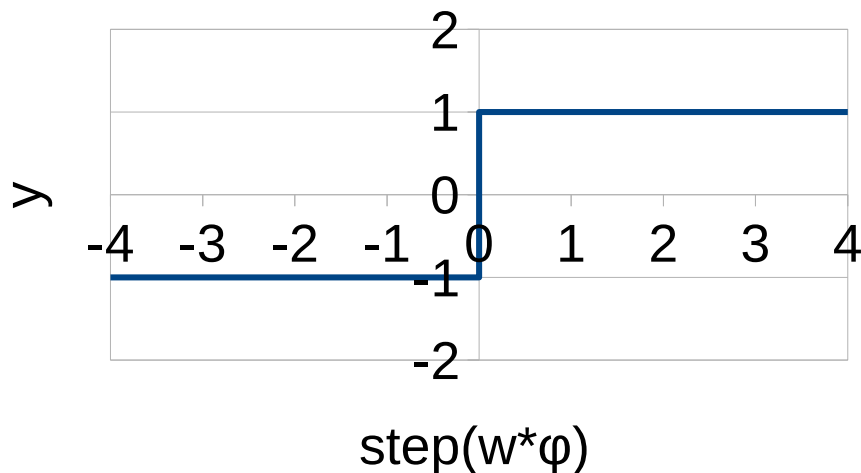
# Example:

- Final neural network:

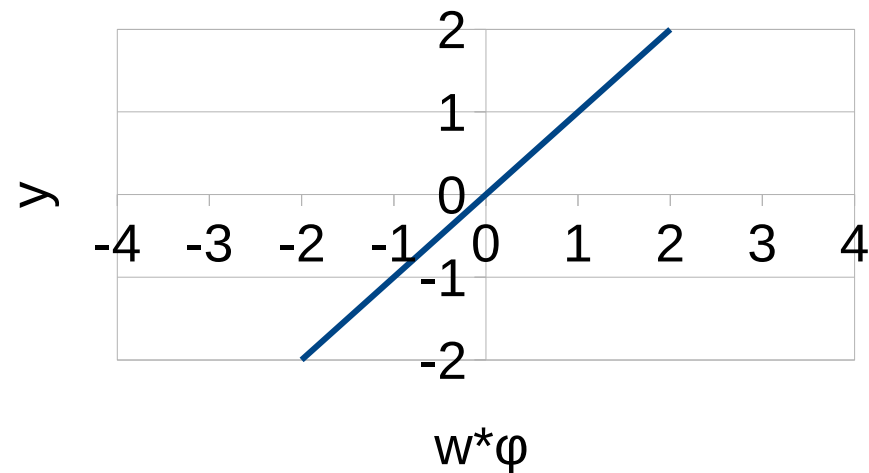# Hidden Layer Activation Functions

## Step Function:
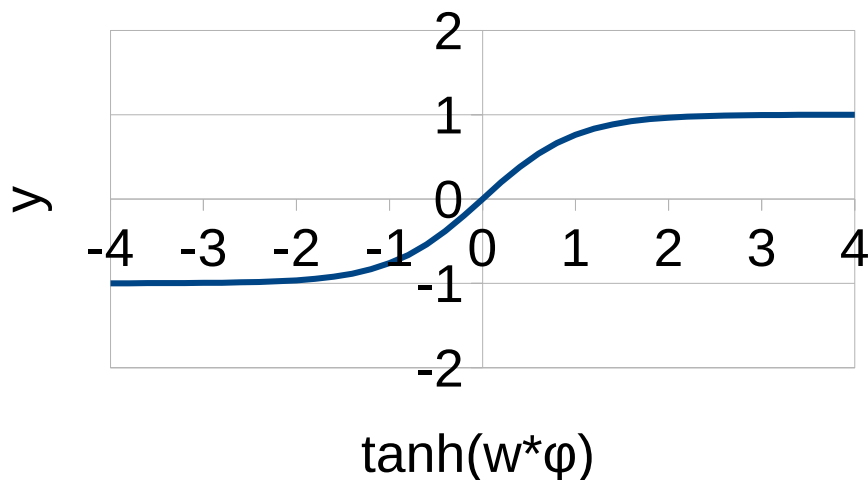- Cannot calculate gradient
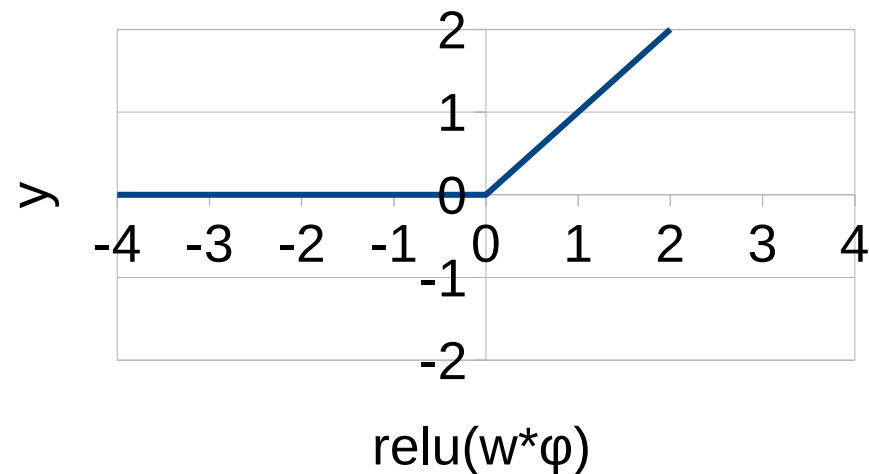
step(w*φ)

## Linear Function:
- Whole net become linear

w*φ

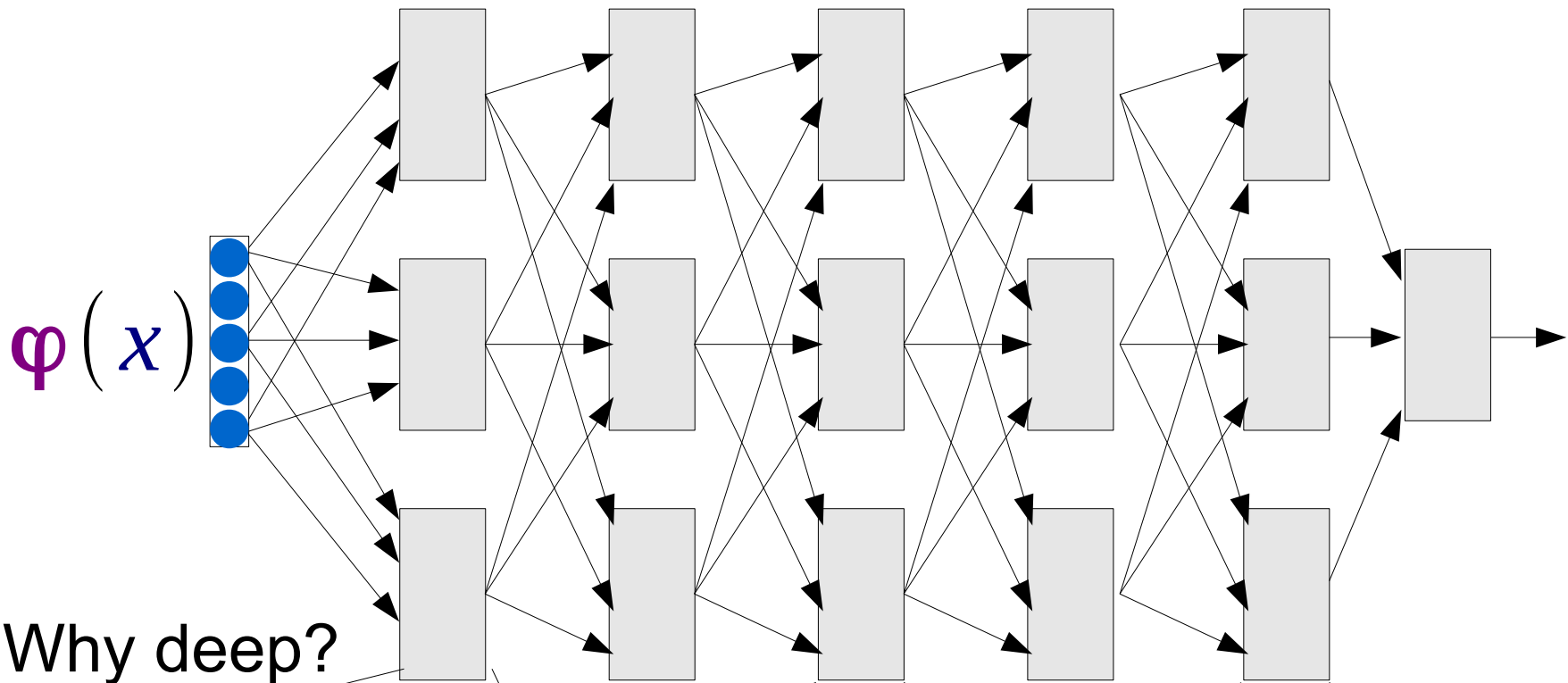## Tanh Function:
Standard (also 0-1 sigmoid)

tanh(w*φ)

## Rectified Linear Function:
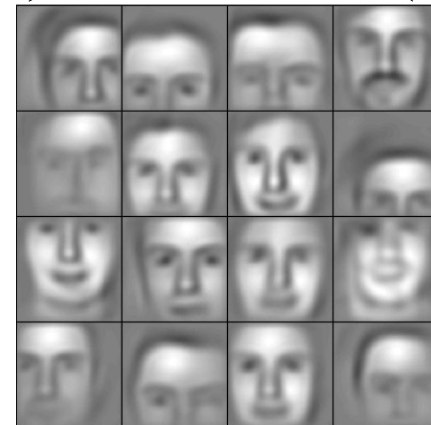+ Gradients at large vals.

relu(w*φ)

18
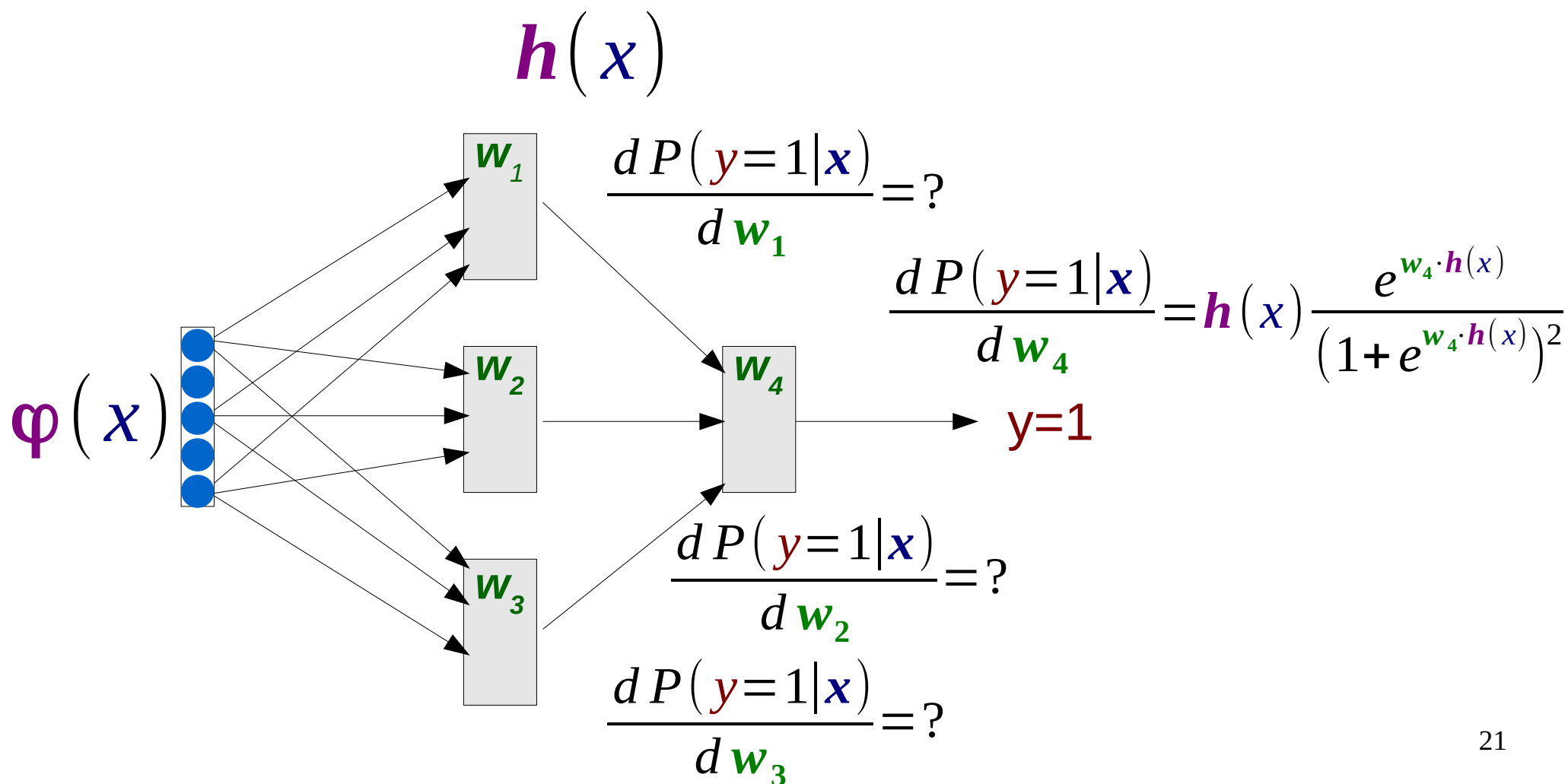
# Deep Networks



$\varphi(x)$

- Why deep?

- Gradually more detailed functions (e.g. [Lee+ 09])

19

# Learning Neural Nets

# Learning:
# Don't Know Derivative for Hidden Units!

- For NNs, only know correct tag for last layer

$$h(x)$$



$$\frac{dP(y=1|x)}{dw_1} = ?$$

$$\frac{dP(y=1|x)}{dw_4} = h(x)\frac{e^{w_4 \cdot h(x)}}{(1+e^{w_4 \cdot h(x)})^2}$$

$$\varphi(x)$$

$$y=1$$

$$\frac{dP(y=1|x)}{dw_2} = ?$$

$$\frac{dP(y=1|x)}{dw_3} = ?$$

21

# Answer: Back-Propogation

- Calculate derivative w/ chain rule

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

$$\frac{e^{w_4 \cdot h(x)}}{(1+e^{w_4 \cdot h(x)})^2}$$
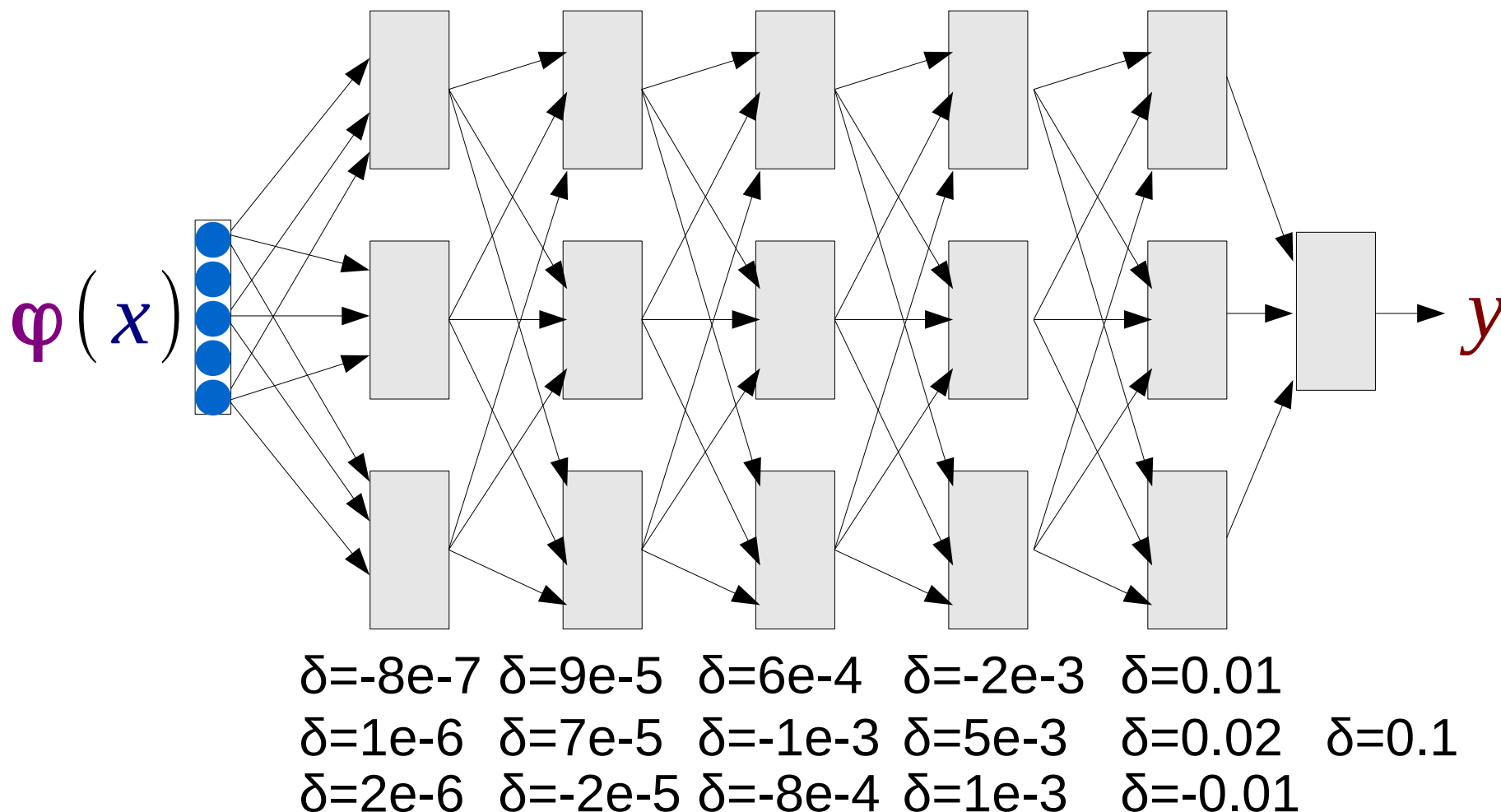
$$W_{1,4}$$

Error of next unit ($\delta_4$)　　Weight　Gradient of this unit

**In General**
Calculate *i* based on next units *j*:
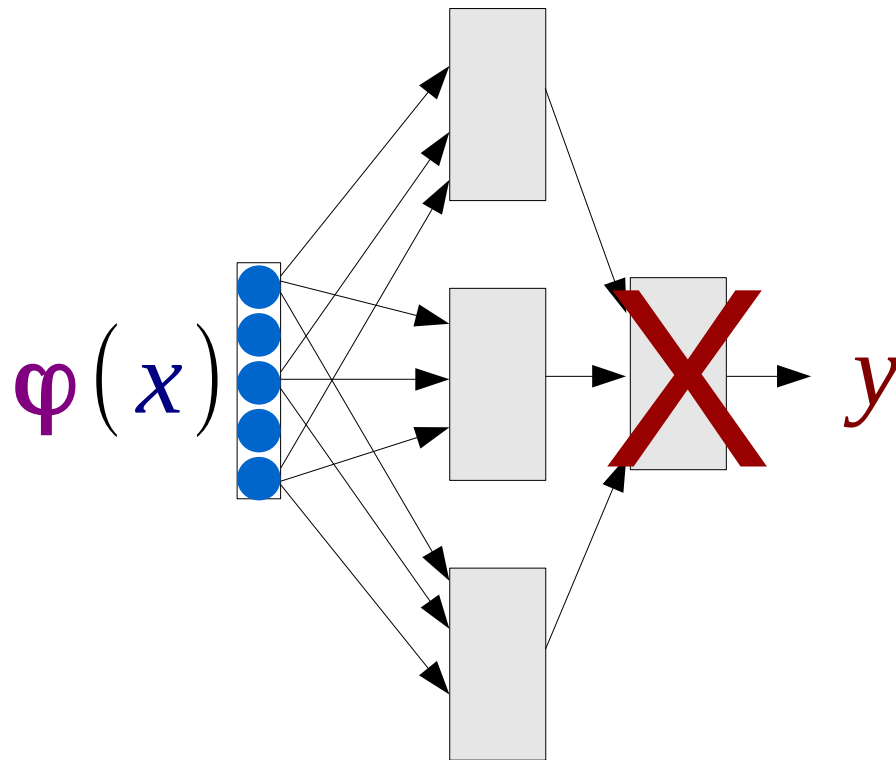
$$\frac{dP(y=1|x)}{w_i} = \frac{dh_i(x)}{dw_i} \sum_j \delta_j w_{i,j}$$
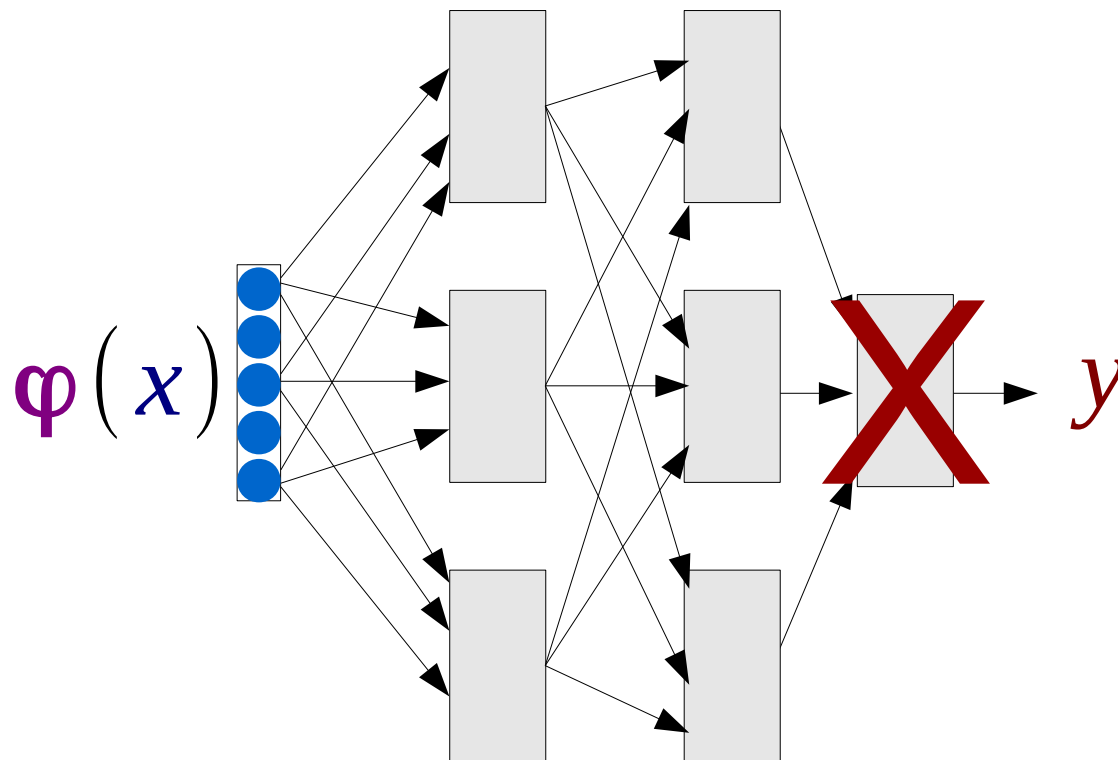
# BP in Deep Networks: Vanishing Gradients



$\varphi(x)$      $y$

$\delta$=-8e-7  $\delta$=9e-5  $\delta$=6e-4  $\delta$=-2e-3  $\delta$=0.01
$\delta$=1e-6  $\delta$=7e-5  $\delta$=-1e-3  $\delta$=5e-3  $\delta$=0.02  $\delta$=0.1
$\delta$=2e-6  $\delta$=-2e-5  $\delta$=-8e-4  $\delta$=1e-3  $\delta$=-0.01

- Exploding gradient as well

# Layerwise Training



- Train one layer at a time

# Layerwise Training
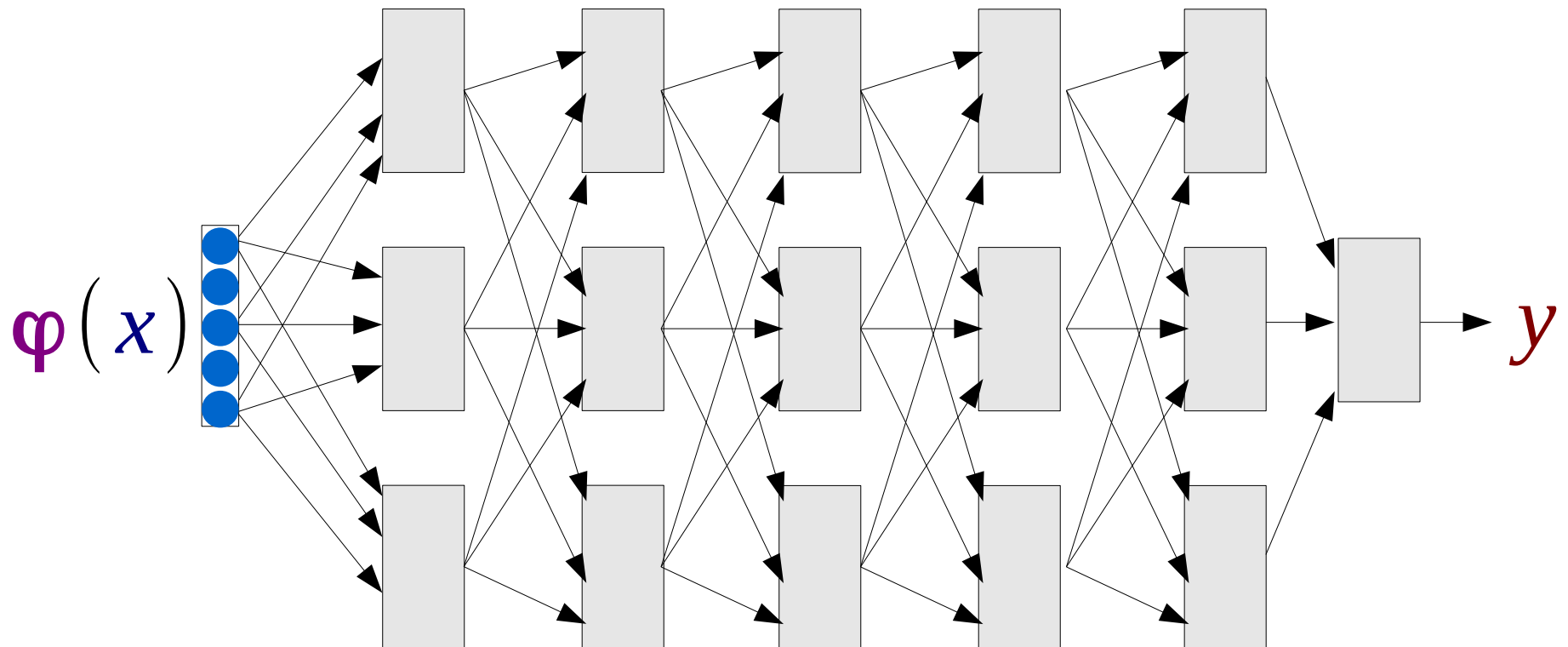


$$\varphi(x) \qquad y$$

- Train one layer at a time

# Layerwise Training



$$\varphi(x)$$

$$y$$

- Train one layer at a time

# Autoencoders

- Initialize the NN by training it to reproduce itself

$$\varphi(x) \qquad\qquad\qquad\qquad\qquad\qquad\qquad \varphi(x)$$
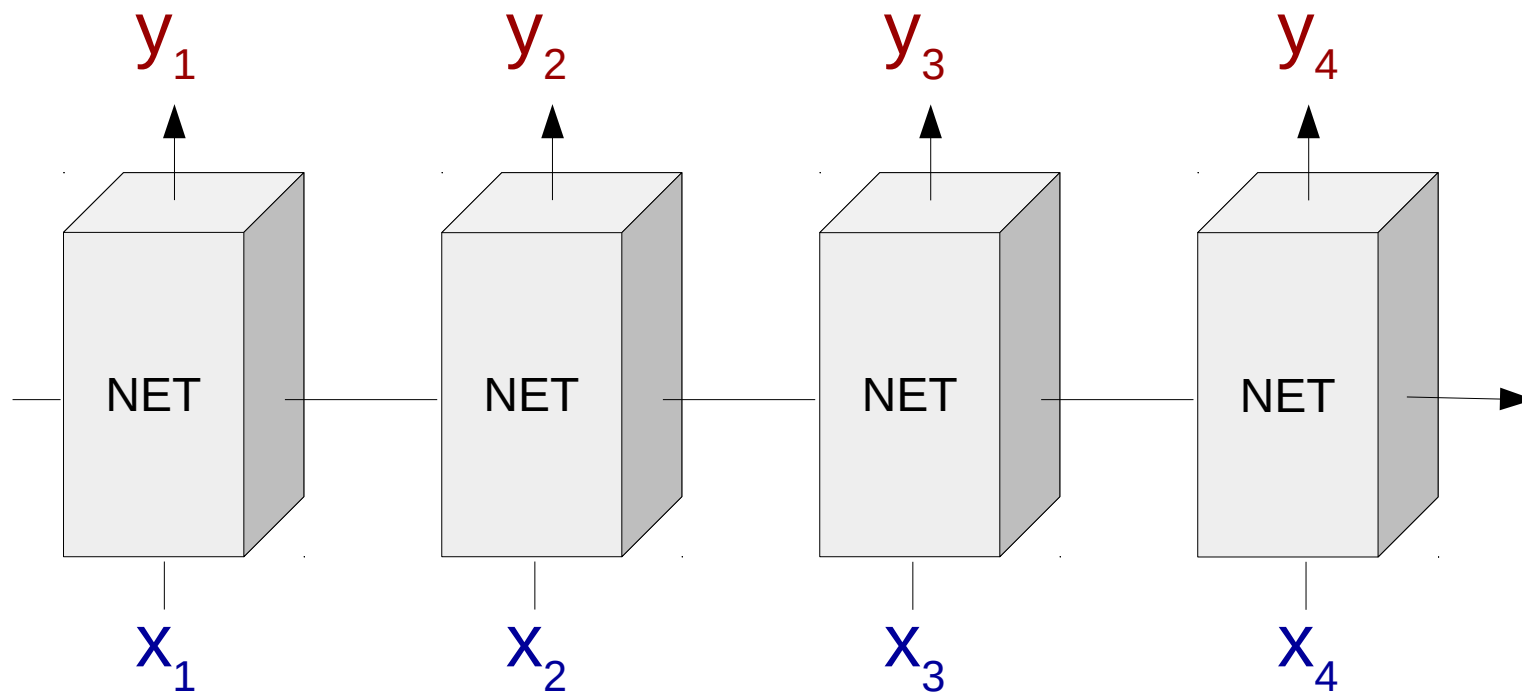
- Advantage: No need for labels y!

# Dropout

- Problem: Overfitting



- Deactivate part of the network on each update
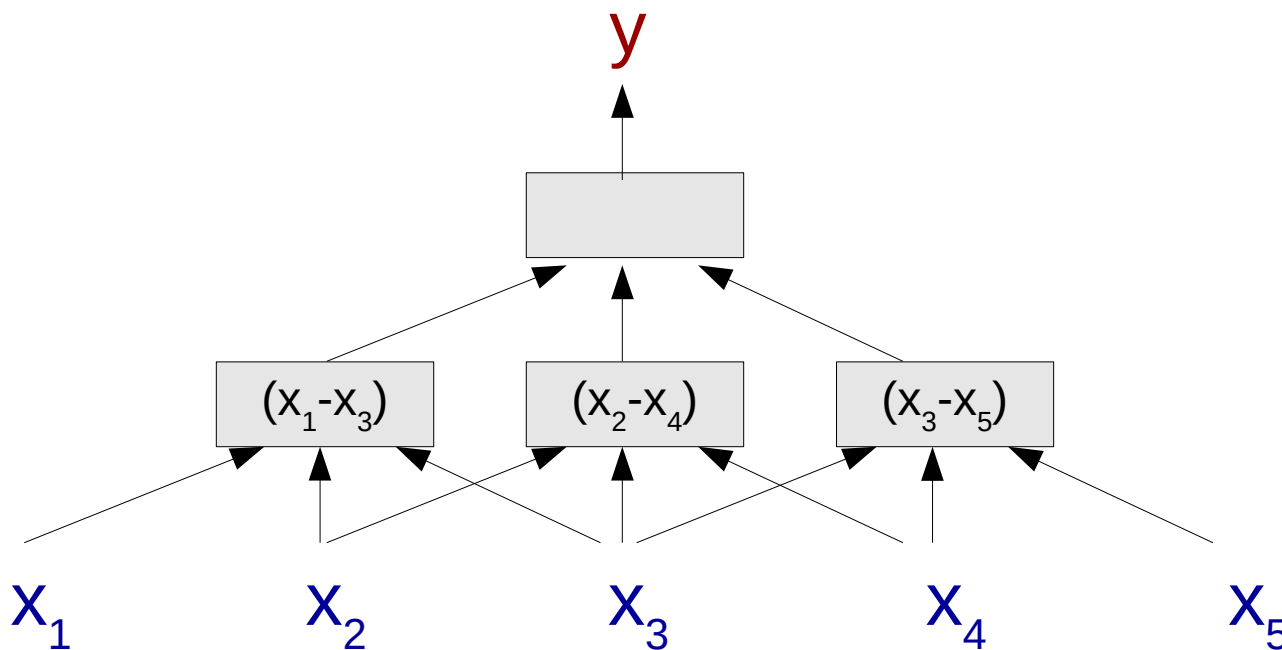- Can be seen as an ensemble method

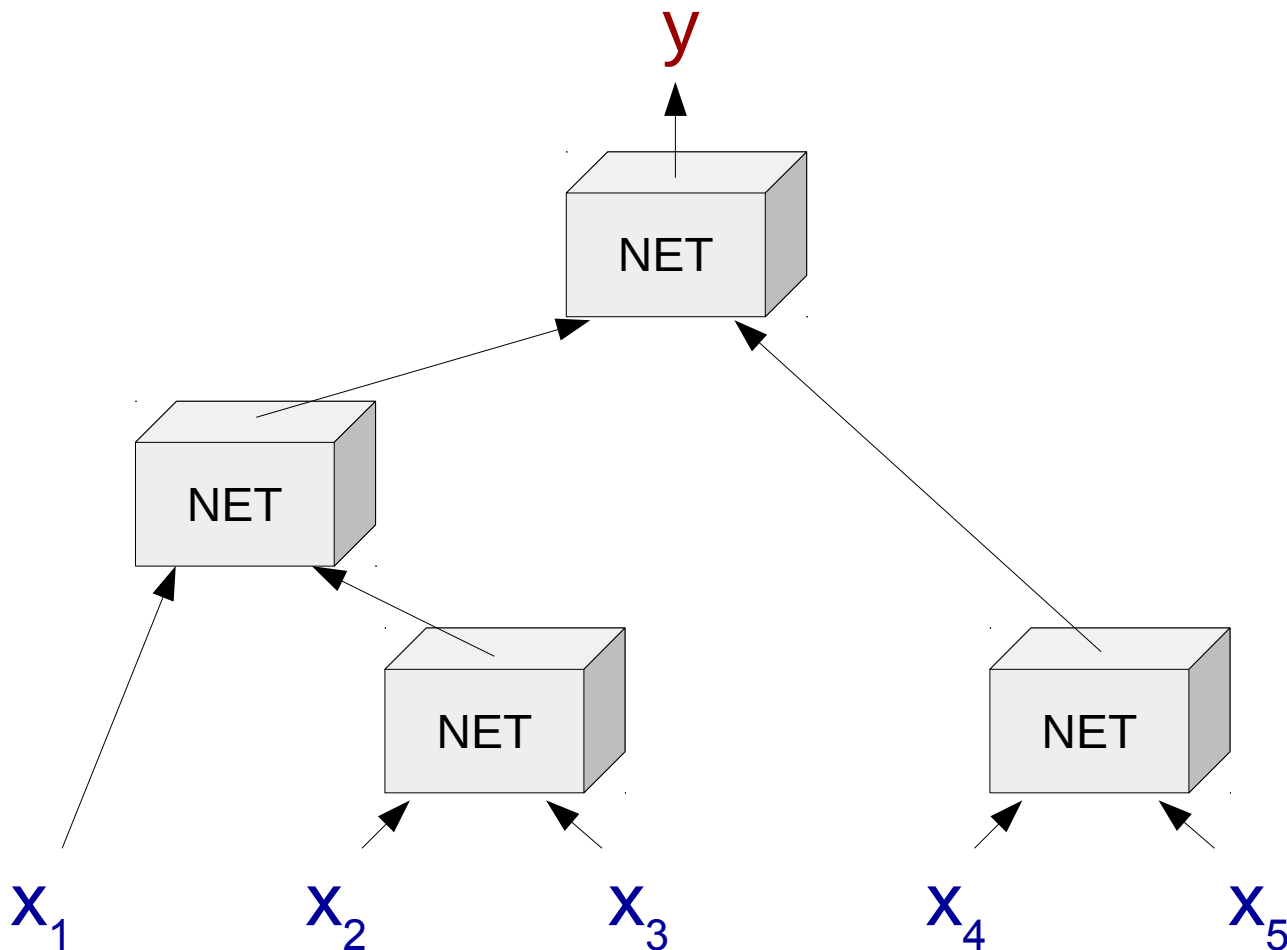# Network Architectures

# Recurrent Neural Nets



- Good for modeling sequence data
- e.g. Speech

# Convolutional Neural Nets

$y$

$(x_1\text{-}x_3)$　$(x_2\text{-}x_4)$　$(x_3\text{-}x_5)$

$x_1$　$x_2$　$x_3$　$x_4$　$x_5$

- Good for modeling data with spatial correlation
- e.g. Images

# Recursive Neural Nets



- Good for modeling data with tree structures
- e.g. Language

# Other Topics

# Other Topics

- Deep belief networks

    - Restricted Boltzmann machines
    - Contrastive estimation

- Generating with Neural Nets

- Visualizing internal states

- Batch/mini-batch update strategies

- Long short-term memory

- Learning on GPUs

- Tools for implementing deep learning