

Logistic Regression & Neural Networks

CMSC 723 / LING 723 / INST 725

Marine Carpuat

Slides credit: Graham Neubig, Jacob
Eisenstein

Logistic Regression

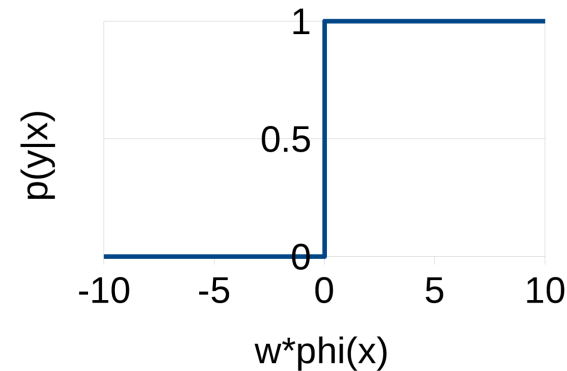
Perceptron & Probabilities

- What if we want a probability $p(y|x)$?
- The perceptron gives us a prediction y
 - Let's illustrate this with binary classification

In other words:

$$P(y=1|x)=1 \text{ if } w \cdot \varphi(x) \geq 0$$

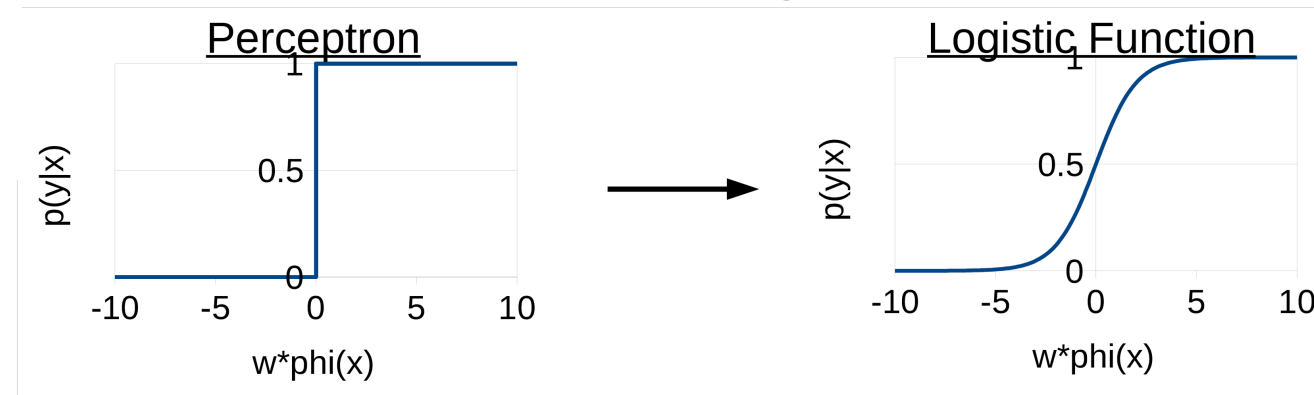
$$P(y=1|x)=0 \text{ if } w \cdot \varphi(x) < 0$$



The logistic function

- x : the input
- $\varphi(x)$: vector of feature functions $\{\varphi_1(x), \varphi_2(x), \dots, \varphi_l(x)\}$
- w : the weight vector $\{w_1, w_2, \dots, w_l\}$
- y : the prediction, +1 if “yes”, -1 if “no”

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



- “Softer” function than in perceptron
- Can account for uncertainty
- Differentiable

Logistic regression: how to train?

- Train based on **conditional likelihood**
- Find parameters \mathbf{w} that maximize conditional likelihood of all answers y_i given examples x_i

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

Stochastic gradient ascent (or descent)

- Online training algorithm for logistic regression
 - and other probabilistic models

```
create map  $w$   
for / iterations  
  for each labeled pair  $x, y$  in the data  
     $w \ += \ \alpha \ * \ dP(y|x)/dw$ 
```

- Update weights for every training example
- Move in direction given by gradient
- Size of update step scaled by learning rate

Gradient of the logistic function

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

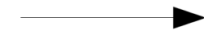
$$\begin{aligned}\frac{d}{d \mathbf{w}} P(\mathbf{y} = -1 | \mathbf{x}) &= \frac{d}{d \mathbf{w}} \left(1 - \frac{e^{\mathbf{w} \cdot \boldsymbol{\varphi}(\mathbf{x})}}{1 + e^{\mathbf{w} \cdot \boldsymbol{\varphi}(\mathbf{x})}} \right) \\ &= -\boldsymbol{\varphi}(\mathbf{x}) \frac{e^{\mathbf{w} \cdot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

Example: Person/not-person classification problem

Given an introductory sentence in Wikipedia
predict whether the article is about a person

Given

Gonso was a Sanron sect priest (754-827)
in the late Nara and early Heian periods.



Predict

Yes!

Shichikuzan Chigogataki Fudomyoo is
a historical site located at Magura, Maizuru
City, Kyoto Prefecture.



No!

Example: initial update

- Set $\alpha=1$, initialize $\mathbf{w}=\mathbf{0}$

\mathbf{x} = A site , located in Maizuru , Kyoto $y = -1$

$$\begin{aligned} \mathbf{w} \cdot \boldsymbol{\varphi}(\mathbf{x}) = 0 \quad \frac{d}{d\mathbf{w}} P(y = -1 | \mathbf{x}) &= -\frac{e^0}{(1+e^0)^2} \boldsymbol{\varphi}(\mathbf{x}) \\ &= -0.25 \boldsymbol{\varphi}(\mathbf{x}) \end{aligned}$$

$$\mathbf{w} \leftarrow \mathbf{w} + -0.25 \boldsymbol{\varphi}(\mathbf{x})$$

$$\begin{aligned} w_{\text{unigram "Maizuru"}} &= -0.25 \\ w_{\text{unigram ","}} &= -0.5 \\ w_{\text{unigram "in"}} &= -0.25 \\ w_{\text{unigram "Kyoto"}} &= -0.25 \end{aligned}$$

$$\begin{aligned} w_{\text{unigram "A"}} &= -0.25 \\ w_{\text{unigram "site"}} &= -0.25 \\ w_{\text{unigram "located"}} &= -0.25 \end{aligned}$$

Example: second update

x = Shoken , monk born in Kyoto

$y = 1$

$$w \cdot \varphi(x) = -1 \quad \begin{array}{c} -0.5 \\ \swarrow \searrow \\ \frac{d}{dw} P(y=1|x) \end{array} \quad \begin{array}{c} -0.25 \quad -0.25 \end{array}$$

$$\frac{d}{dw} P(y=1|x) = \frac{e^1}{(1+e^1)^2} \varphi(x)$$

$$= 0.196 \varphi(x)$$

$$w \leftarrow w + 0.196 \varphi(x)$$

$w_{\text{unigram "Maizuru"}} = -0.25$
 $w_{\text{unigram ","}} = -0.304$
 $w_{\text{unigram "in"}} = -0.054$
 $w_{\text{unigram "Kyoto"}} = -0.054$

$w_{\text{unigram "A"}} = -0.25$
 $w_{\text{unigram "site"}} = -0.25$
 $w_{\text{unigram "located"}} = -0.25$

$w_{\text{unigram "Shoken"}} = 0.196$
 $w_{\text{unigram "monk"}} = 0.196$
 $w_{\text{unigram "born"}} = 0.196$

How to set the learning rate?

- Various strategies
 - decay over time

$$\alpha = \frac{1}{C + t}$$

Parameter

Number of
samples

- Use held-out test set, increase learning rate when likelihood increases

Multiclass version

$$p(y \mid \boldsymbol{x}) = \frac{\exp(\boldsymbol{\theta}^\top \boldsymbol{f}(\boldsymbol{x}, y))}{\sum_{y' \in \mathcal{Y}} \exp(\boldsymbol{\theta}^\top \boldsymbol{f}(\boldsymbol{x}, y'))}.$$

Some models are better than others...

- Consider these 2 examples

-1 he saw a bird in the park
+1 he saw a robbery in the park

- Which of the 2 models below is better?

Classifier 1

he +3
saw -5
a +0.5
bird -1
robbery +1
in +5
the -3
park -2

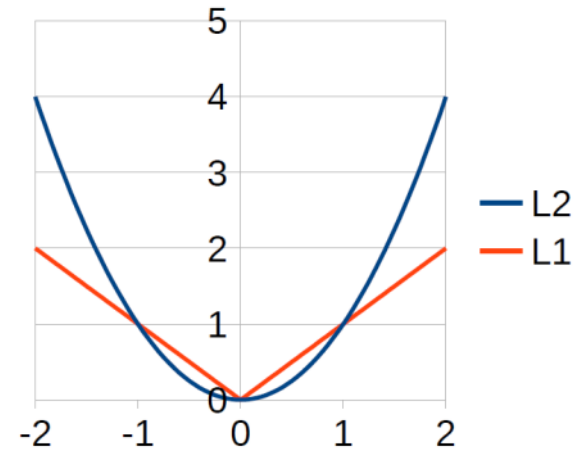
Classifier 2

bird -1
robbery +1

Classifier 2 will probably generalize better!
It does not include irrelevant information
=> Smaller model is better

Regularization

- A penalty on adding extra weights
- L2 regularization: $\|w\|_2$
 - big penalty on large weights
 - small penalty on small weights
- L1 regularization: $\|w\|_1$
 - Uniform increase when large or small
 - Will cause many weights to become zero



L1 regularization in online learning

```
update_weights(w, phi, y, c)
★ for name, value in w:
★     if abs(value) < c:
★         w[name] = 0
★     else:
★         w[name] -= sign(value) * c
★ for name, value in phi:
    w[name] += value * y
```

If abs. value < *c*,
set weight to zero

If value > 0,
decrease by *c*

If value < 0,
increase by *c*

What you should know

- Standard supervised learning set-up for text classification
 - Difference between train vs. test data
 - How to evaluate
- 3 examples of supervised linear classifiers
 - Naïve Bayes, Perceptron, Logistic Regression
 - Learning as optimization: what is the objective function optimized?
 - Difference between generative vs. discriminative classifiers
 - Smoothing, regularization
 - Overfitting, underfitting

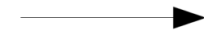
Neural networks

Person/not-person classification problem

Given an introductory sentence in Wikipedia
predict whether the article is about a person

Given

Gonso was a Sanron sect priest (754-827)
in the late Nara and early Heian periods.



Predict

Yes!

Shichikuzan Chigogataki Fudomyoo is
a historical site located at Magura, Maizuru
City, Kyoto Prefecture.



No!

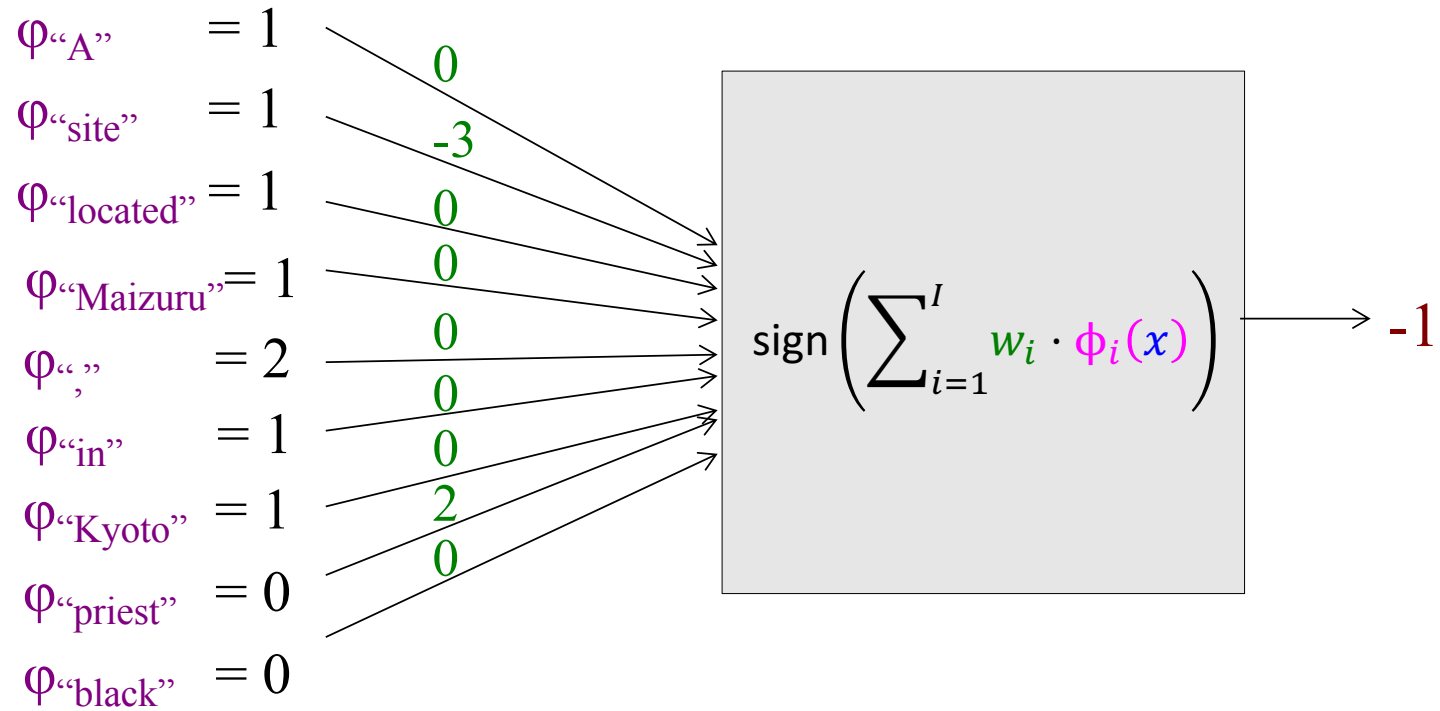
Formalizing binary prediction

$$\begin{aligned} y &= \text{sign}(\mathbf{w} \cdot \boldsymbol{\varphi}(\mathbf{x})) \\ &= \text{sign}\left(\sum_{i=1}^I \mathbf{w}_i \cdot \varphi_i(\mathbf{x})\right) \end{aligned}$$

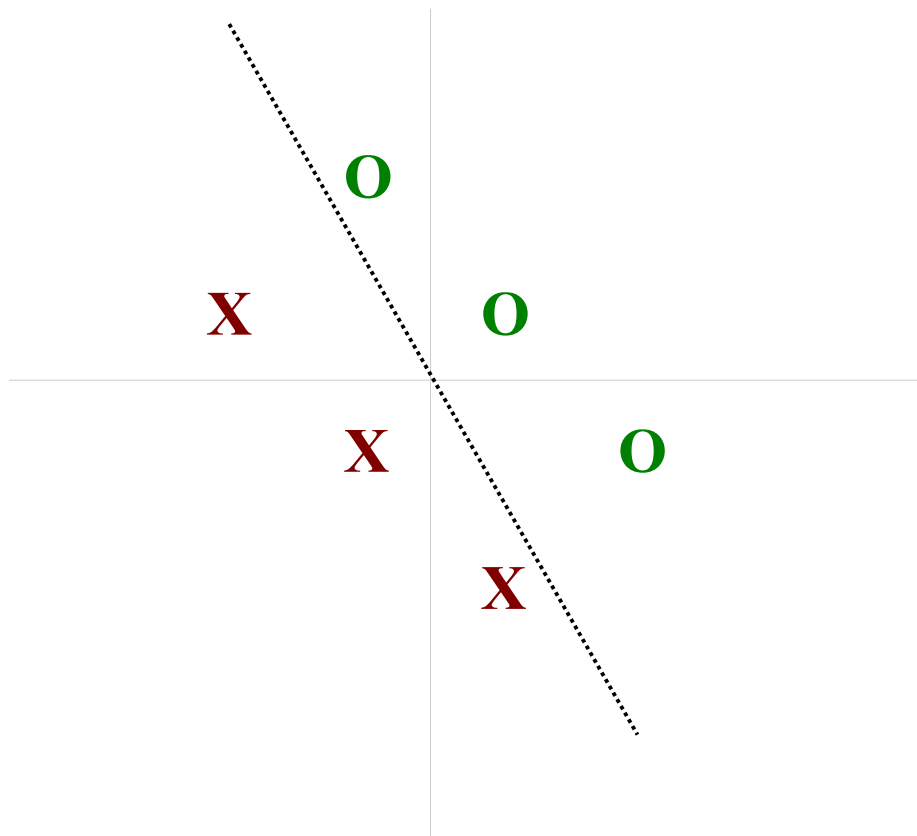
- \mathbf{x} : the input
- $\boldsymbol{\varphi}(\mathbf{x})$: vector of feature functions $\{\varphi_1(\mathbf{x}), \varphi_2(\mathbf{x}), \dots, \varphi_I(\mathbf{x})\}$
- \mathbf{w} : the weight vector $\{w_1, w_2, \dots, w_I\}$
- y : the prediction, +1 if “yes”, -1 if “no”
 - ($\text{sign}(v)$ is +1 if $v \geq 0$, -1 otherwise)

The Perceptron:

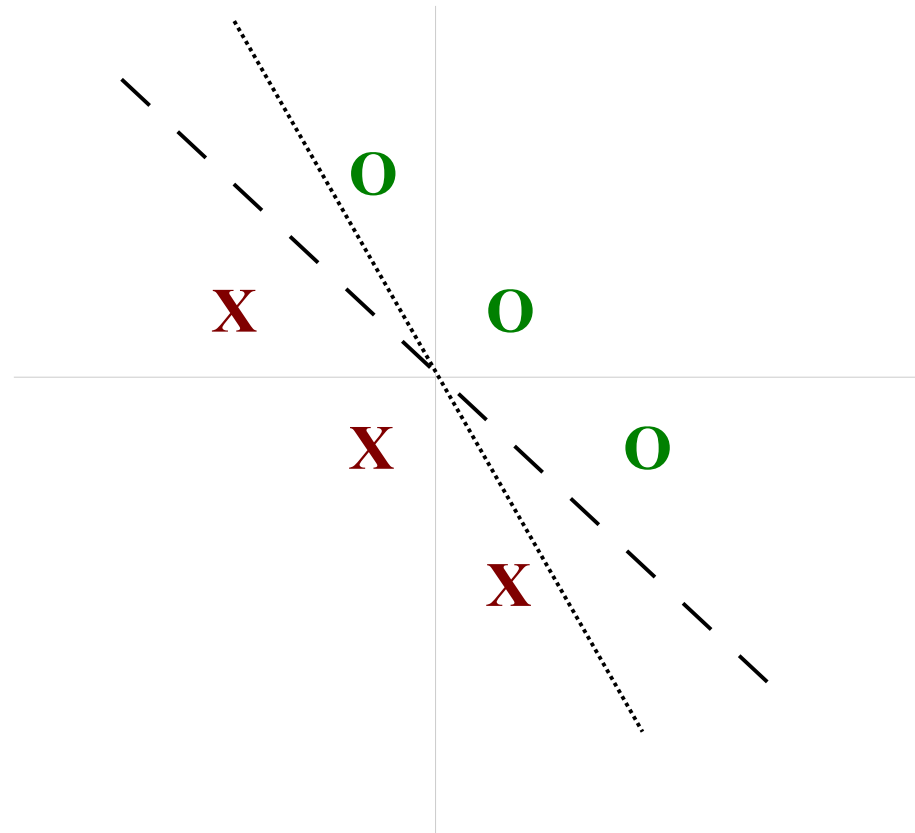
a “machine” to calculate a weighted sum



The Perceptron: Geometric interpretation

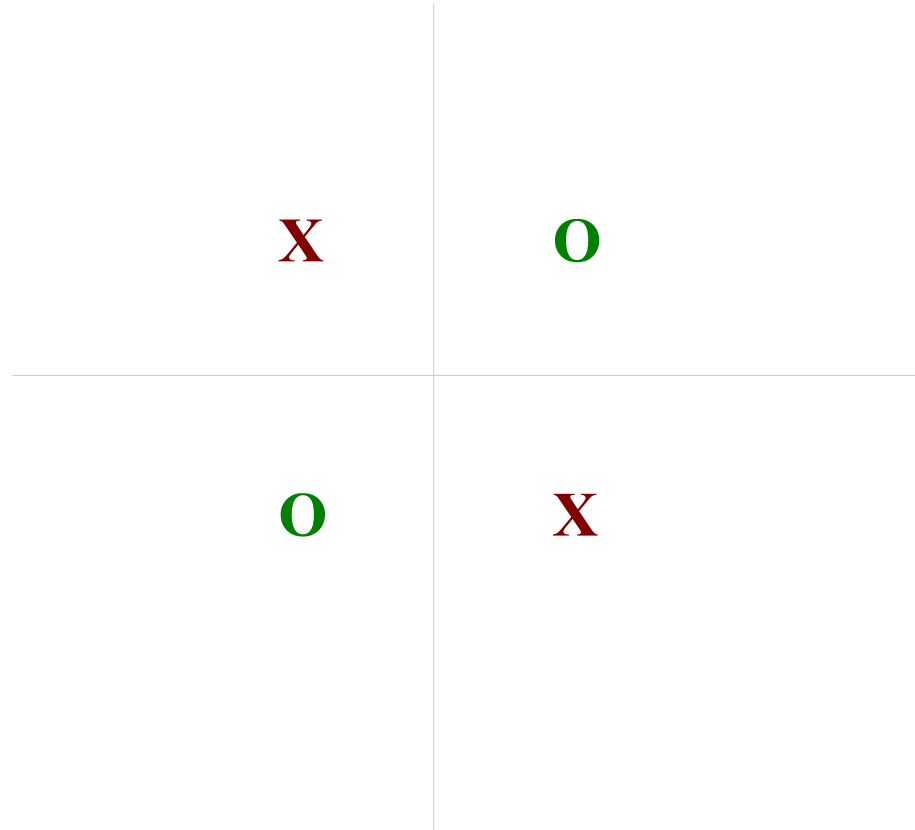


The Perceptron: Geometric interpretation



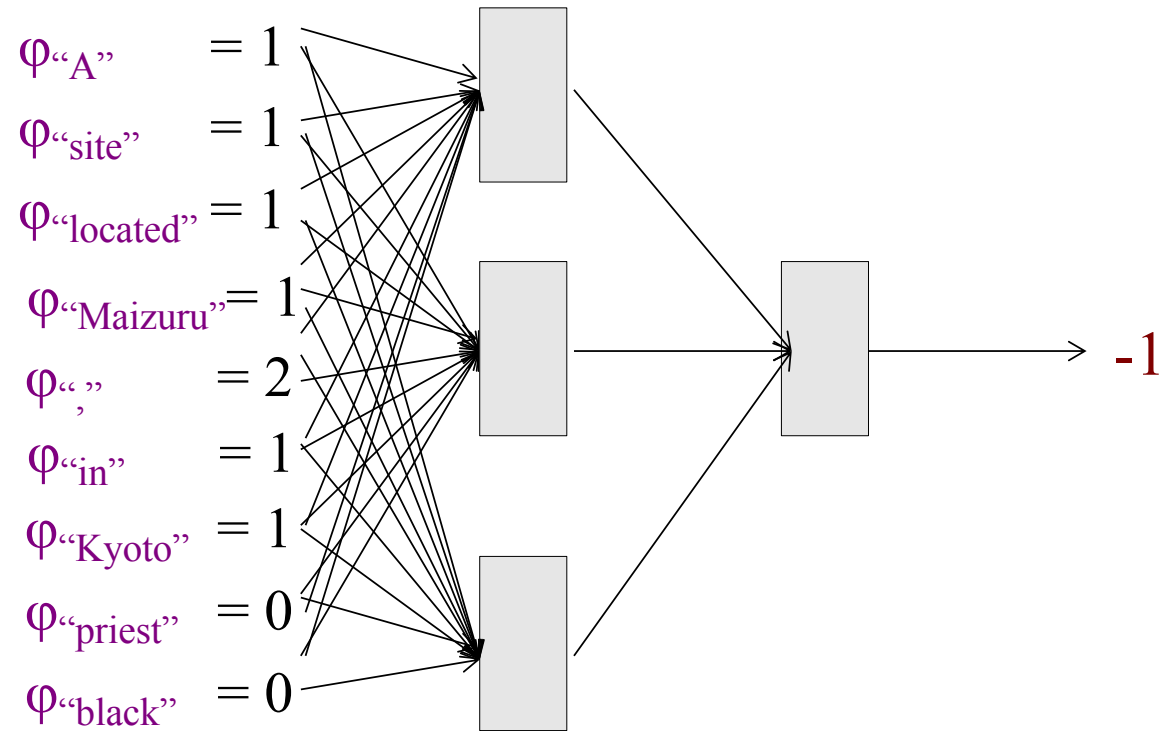
Limitation of perceptron

- can only find **linear separations** between positive and negative examples



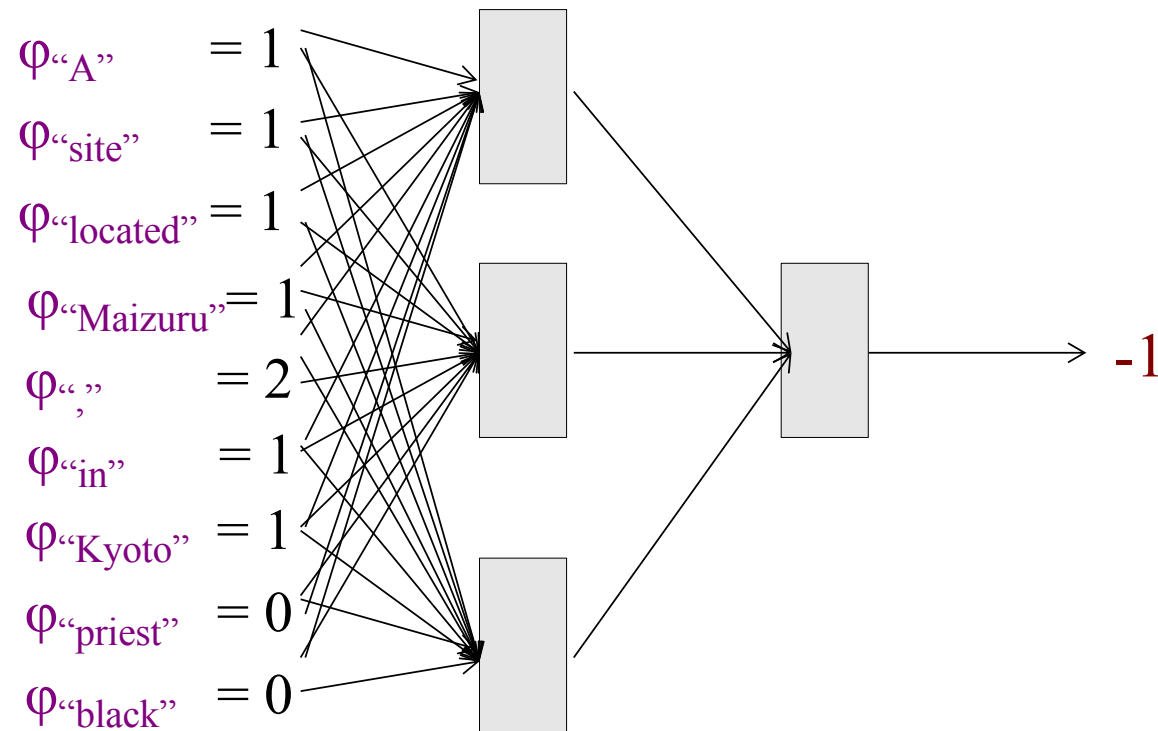
Neural Networks

- Connect together multiple perceptrons



- Motivation: Can represent non-linear functions!

Neural Networks: key terms

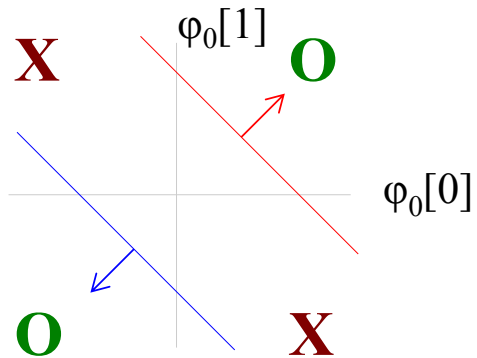


- Input (aka features)
- Output
- Nodes
- Layers
- Hidden layers
- Activation function (non-linear)
- Multi-layer perceptron

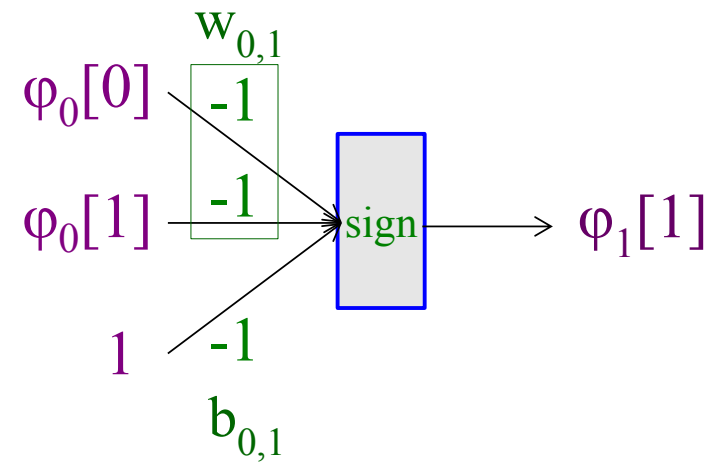
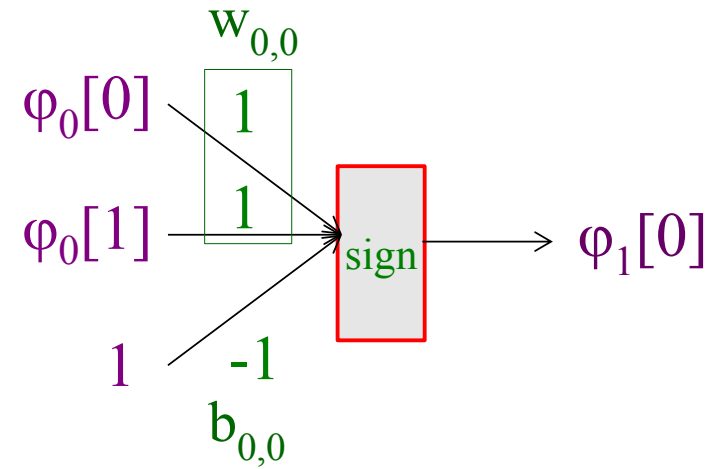
Example

- Create two classifiers

$$\varphi_0(\mathbf{x}_1) = \{-1, 1\} \quad \varphi_0(\mathbf{x}_2) = \{1, 1\}$$

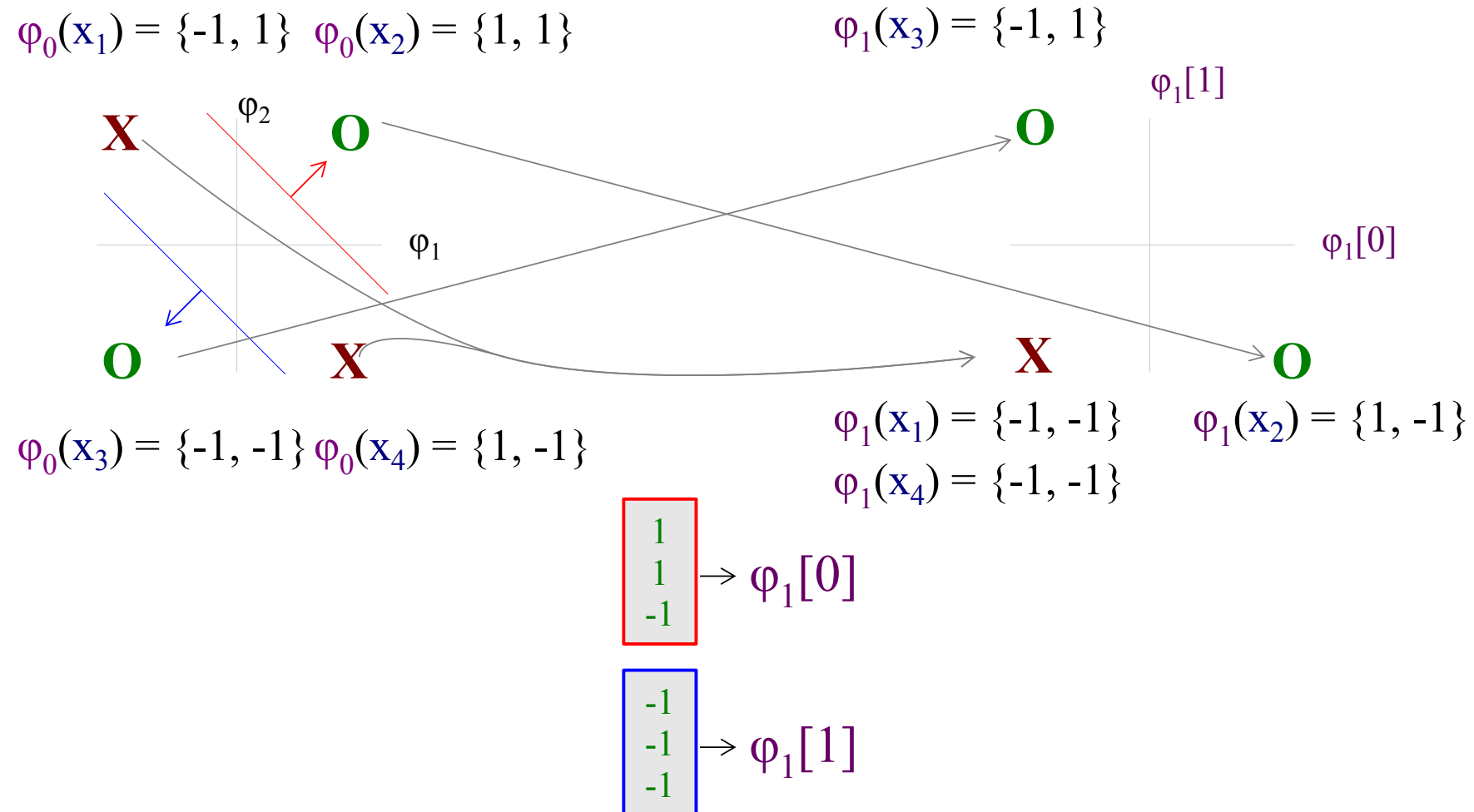


$$\varphi_0(\mathbf{x}_3) = \{-1, -1\} \quad \varphi_0(\mathbf{x}_4) = \{1, -1\}$$



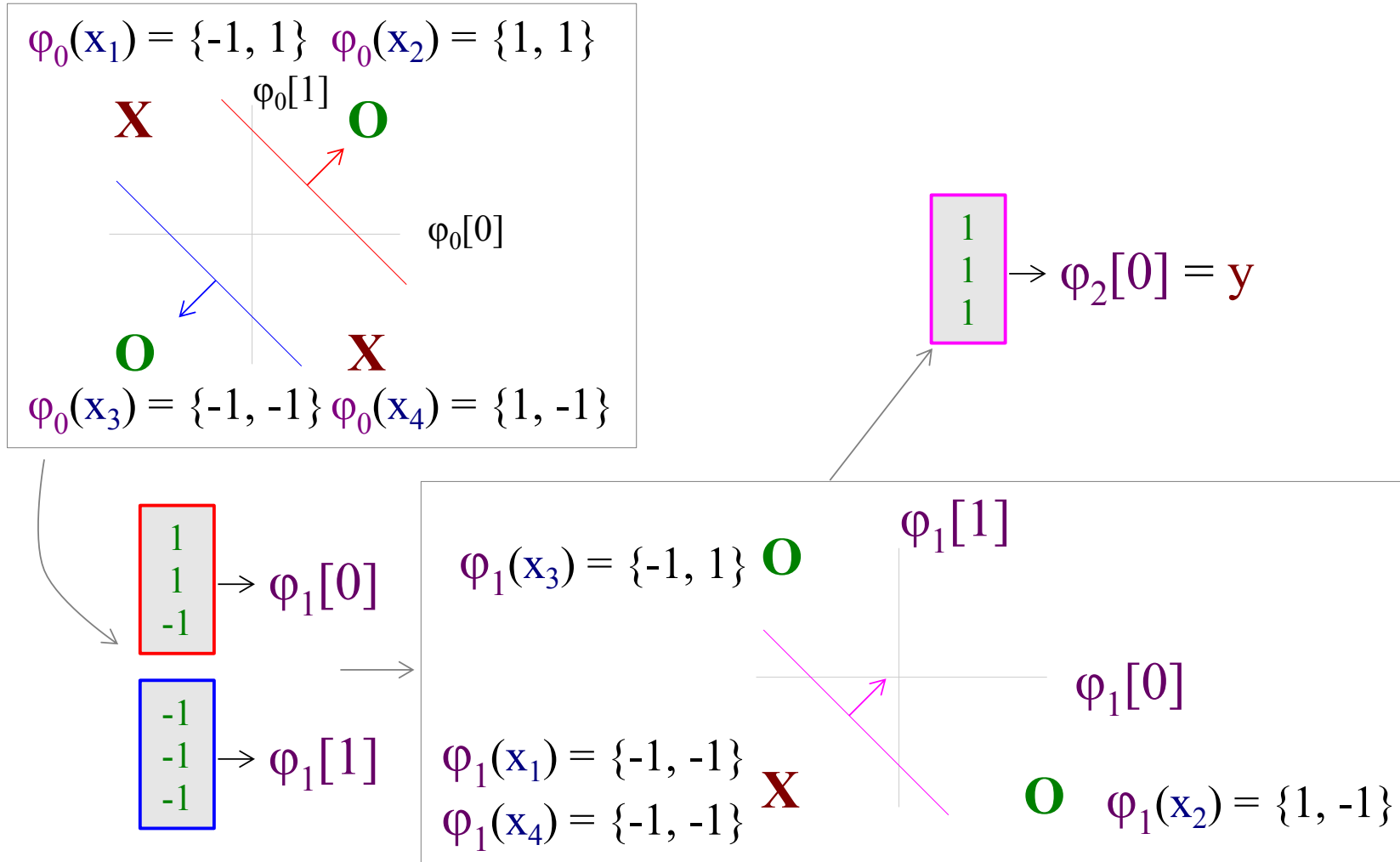
Example

- These classifiers map to a new space



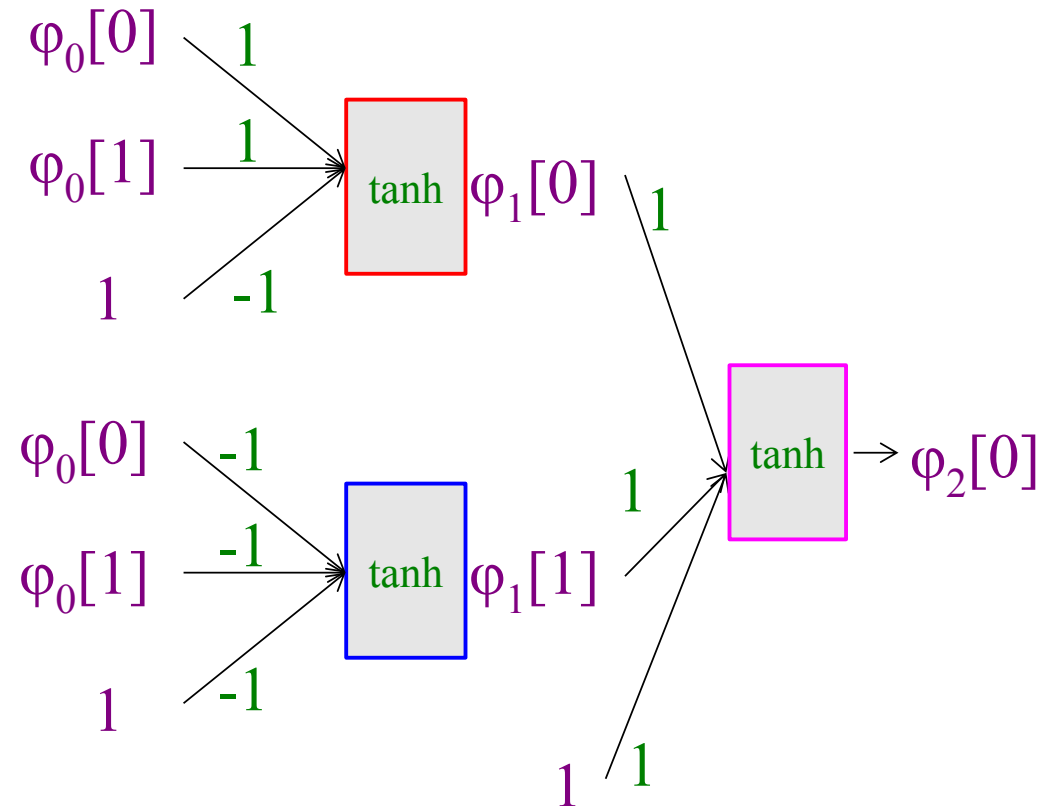
Example

- In new space, the examples are linearly separable!



Example wrap-up: Forward propagation

- The final net



Softmax Function for multiclass classification

- Sigmoid function for multiple classes

$$P(\textcolor{red}{y} \mid \textcolor{blue}{x}) = \frac{e^{\textcolor{green}{\mathbf{w}} \cdot \textcolor{violet}{\phi}(\textcolor{blue}{x}, \textcolor{red}{y})}}{\sum_{\tilde{y}} e^{\textcolor{green}{\mathbf{w}} \cdot \textcolor{violet}{\phi}(\textcolor{blue}{x}, \tilde{y})}}$$

← Current class

← Sum of other classes

- Can be expressed using matrix/vector ops

$$\mathbf{r} = \exp(\textcolor{green}{\mathbf{W}} \cdot \textcolor{violet}{\phi}(\textcolor{blue}{x}, \textcolor{red}{y}))$$

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

Stochastic Gradient Descent

Online training algorithm for probabilistic models

```
 $w = 0$   
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

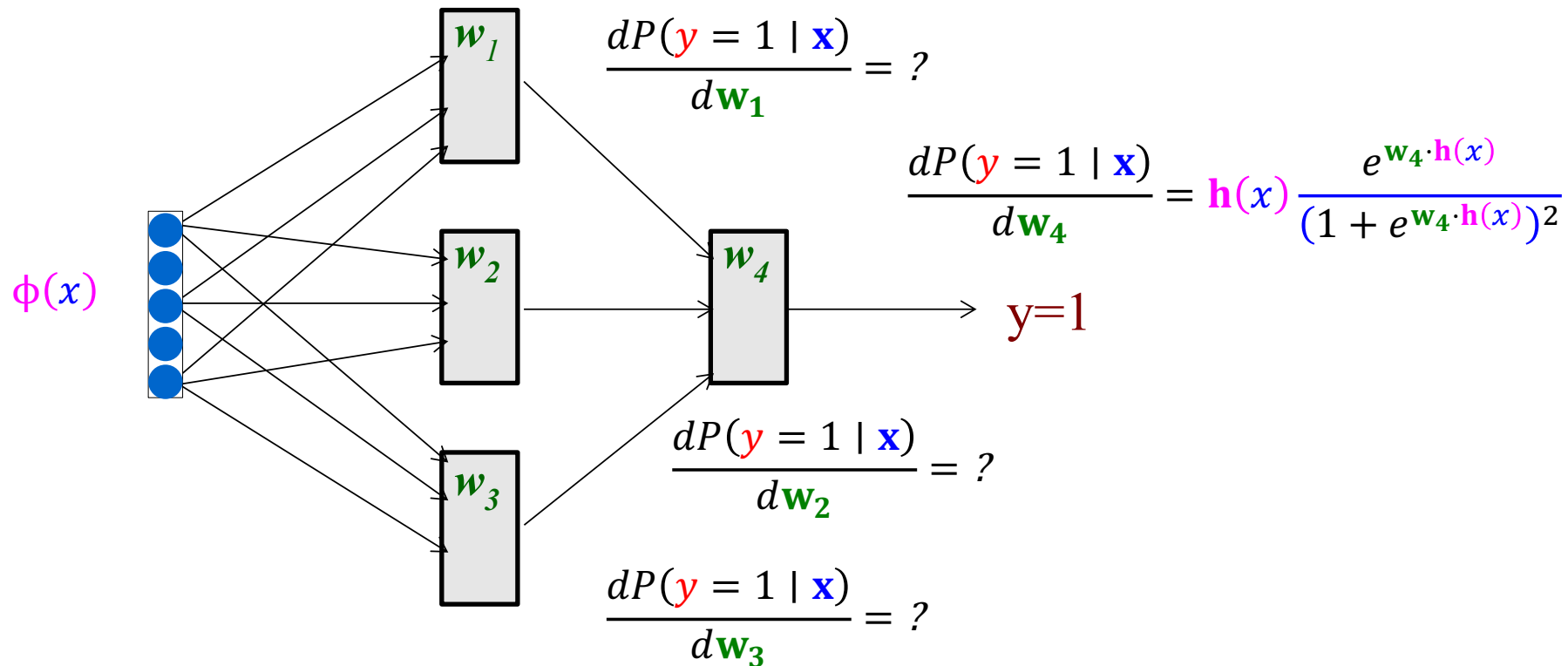
$$\begin{aligned}\frac{d}{d\mathbf{w}} P(\mathbf{y} = 1 \mid \mathbf{x}) &= \frac{d}{d\mathbf{w}} \frac{e^{\mathbf{w} \cdot \phi(\mathbf{x})}}{1 + e^{\mathbf{w} \cdot \phi(\mathbf{x})}} \\ &= \phi(\mathbf{x}) \frac{e^{\mathbf{w} \cdot \phi(\mathbf{x})}}{(1 + e^{\mathbf{w} \cdot \phi(\mathbf{x})})^2}\end{aligned}$$

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

Learning: We Don't Know the Derivative for Hidden Units!

For NNs, only know correct tag for last layer

$\mathbf{h}(x)$



Answer: Back-Propagation

Calculate derivative with chain rule

$$\frac{dP(\textcolor{red}{y} = 1 \mid \textcolor{blue}{x})}{d\textcolor{green}{w}_1} = \frac{dP(\textcolor{red}{y} = 1 \mid \textcolor{blue}{x})}{d\textcolor{green}{w}_4 \textcolor{magenta}{h}(\textcolor{blue}{x})} \frac{d\textcolor{green}{w}_4 \textcolor{magenta}{h}(\textcolor{blue}{x})}{d\textcolor{magenta}{h}_1(\textcolor{blue}{x})} \frac{dh_1(\textcolor{blue}{x})}{d\textcolor{green}{w}_1}$$

$\swarrow \qquad \searrow$

$$\frac{e^{\textcolor{green}{w}_4 \cdot \textcolor{magenta}{h}(\textcolor{blue}{x})}}{(1 + e^{\textcolor{green}{w}_4 \cdot \textcolor{magenta}{h}(\textcolor{blue}{x})})^2} \qquad \textcolor{green}{w}_{1,4}$$

$\downarrow \qquad \downarrow \qquad \downarrow$

Error of
next unit (δ_4) Weight Gradient of
this unit

In General

Calculate i based
on next units j :

$$\frac{dP(\textcolor{red}{y} = 1 \mid \textcolor{blue}{x})}{d\textcolor{green}{w}_i} = \frac{d\textcolor{magenta}{h}_i(\textcolor{blue}{x})}{d\textcolor{green}{w}_i} \sum_j \delta_j \textcolor{green}{w}_{i,j}$$

Backpropagation

=

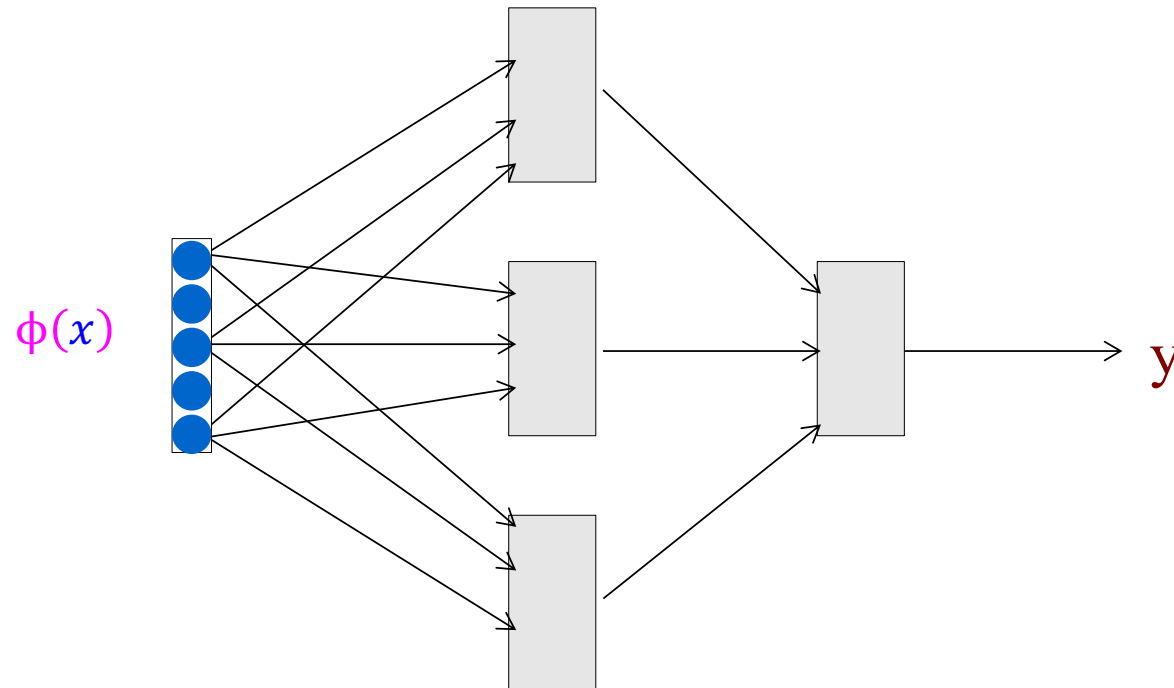
Gradient descent

+

Chain rule

Feed Forward Neural Nets

All connections point **forward**



It is a directed acyclic graph (DAG)

Neural Networks

- Non-linear classification
- Prediction: forward propagation
 - Vector/matrix operations + non-linearities
- Training: backpropagation + stochastic gradient descent

For more details, see [CIML Chap 7](#)