# Machine Learning 2

DS 4420 - Spring 2020

## Neural Networks & backprop

Byron C Wallace

# Neural Networks!

- In 2020, neural networks are the dominant technology in machine learning (for better or worse)!

# Neural Networks!

- In 2020, neural networks are the dominant technology in machine learning (for better or worse)!

- Today, we'll go over some of the fundamentals of NNs and modern libraries (we saw a preview last week, with auto-diff)!

# Neural Networks!

- In 2020, neural networks are the dominant technology in machine learning (for better or worse)!

- Today, we'll go over some of the fundamentals of NNs and modern libraries (we saw a preview last week, with auto-diff)!

- This will also serve as a refresher on gradient descent

# Gradient Descent in Linear Models

Last time we thought in **probabilistic terms** and discussed **maximum likelihood estimation** for "generative" models

# Gradient Descent in Linear Models

Last time we thought in **probabilistic terms** and discussed **maximum likelihood estimation** for "generative" models

Today we'll take the view of learning as **search/optimization**

# Gradient Descent in Linear Models

Last time we thought in **probabilistic terms** and discussed **maximum likelihood estimation** for "generative" models

Today we'll take the view of learning as **search/optimization**

We'll start with linear models, review **gradient descent**, and then talk about **neural nets** + **backprop**

# Loss

The simplest loss is probably 0/1 loss:

<span style="color:green">0</span> if we're correct
<span style="color:red">1</span> if we're wrong

What's an algo that minimizes this?

The *Perceptron*!

Training data $\langle x, y \rangle$

Consider a simple linear model with parameters *w*

$$\hat{y}_i = \begin{cases} 1 & \text{if } \boxed{w \cdot x_i > 0} \\ -1 & \text{otherwise} \end{cases}$$

Training data $\langle x, y \rangle$

Consider a simple linear model with parameters *w*

$$\hat{y}_i = \begin{cases} 1 & \text{if } \boxed{w \cdot x_i > 0} \\ -1 & \text{otherwise} \end{cases}$$

(assumes bias term moved into *x* or omitted)

Training data $\langle x, y \rangle$

Consider a simple linear model with parameters $w$

$$\hat{y}_i = \begin{cases} 1 & \text{if } w \cdot x_i > 0 \\ -1 & \text{otherwise} \end{cases}$$

(assumes bias term moved into $x$ or omitted)

The learning problem is to estimate $w$

Training data $\langle x, y \rangle$

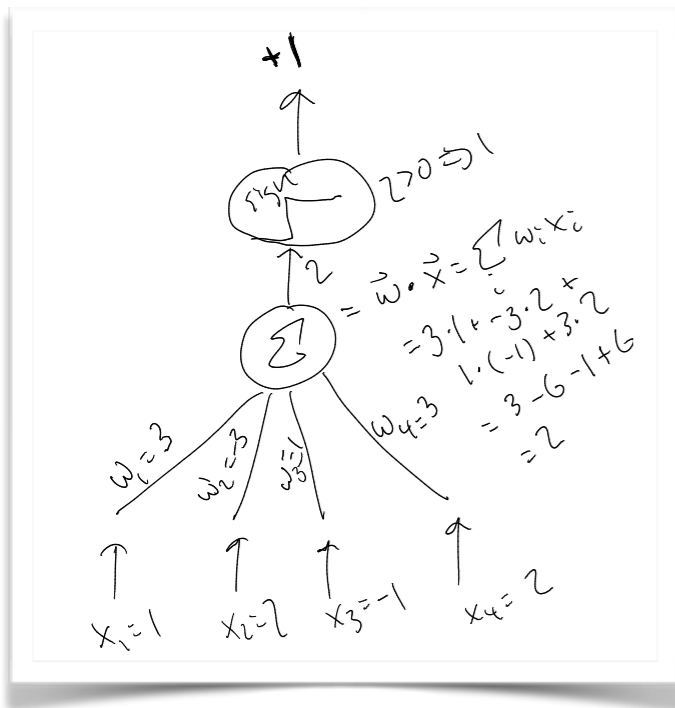Consider a simple linear model with parameters $w$

$$\hat{y}_i = \begin{cases} 1 & \text{if } \boxed{w \cdot x_i > 0} \\ -1 & \text{otherwise} \end{cases}$$

(assumes bias term moved into $x$ or omitted)
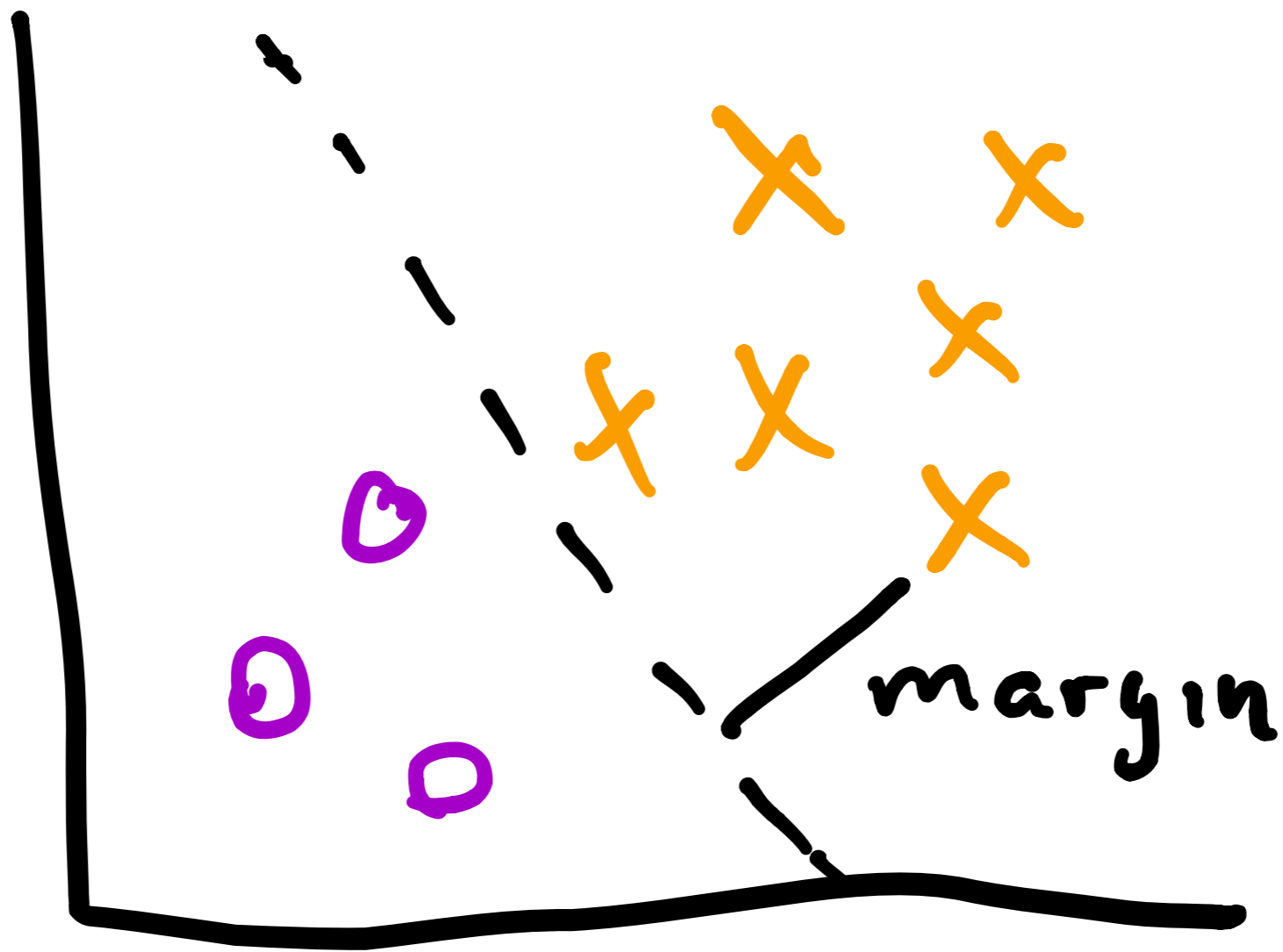
The learning problem is to estimate $w$

What is our criterion for a good $w$?  Minimal **loss**
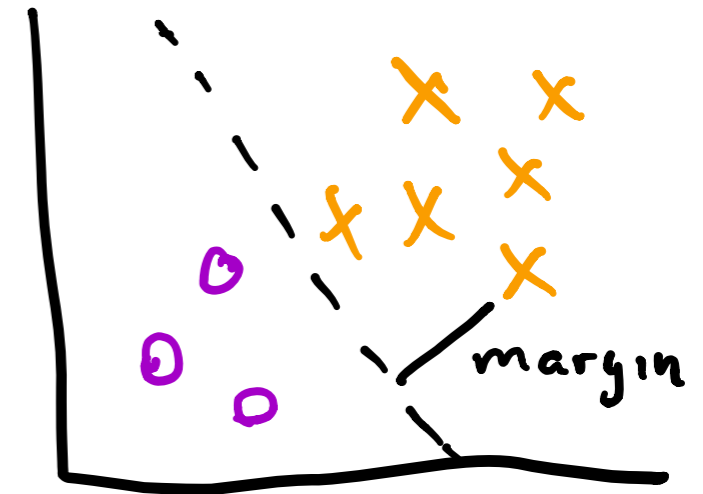
# Perceptron!





**Algorithm 5** PERCEPTRONTRAIN($\mathbf{D}$, *MaxIter*)

1:   $w_d \leftarrow 0$, for all $d = 1 \dots D$            // initialize weights
2:   $b \leftarrow 0$            // initialize bias
3:   **for** $iter = 1 \dots MaxIter$ **do**
4:      **for all** $(x,y) \in \mathbf{D}$ **do**
5:        $a \leftarrow \sum_{d=1}^{D} w_d \, x_d + b$            // compute activation for this example
6:        **if** $ya \leq 0$ **then**
7:          $w_d \leftarrow w_d + yx_d$, for all $d = 1 \dots D$            // update weights
8:          $b \leftarrow b + y$            // update bias
9:        **end if**
10:      **end for**
11:   **end for**
12:   **return** $w_0, w_1, \dots, w_D, b$

*Fig and Alg from CIML [Daume]*

margin

# Problems with 0/1 loss

- If we're wrong by .0001 it is "as bad" as being wrong by .9999

- Because it is discrete, optimization is hard if the instances are not linearly separable
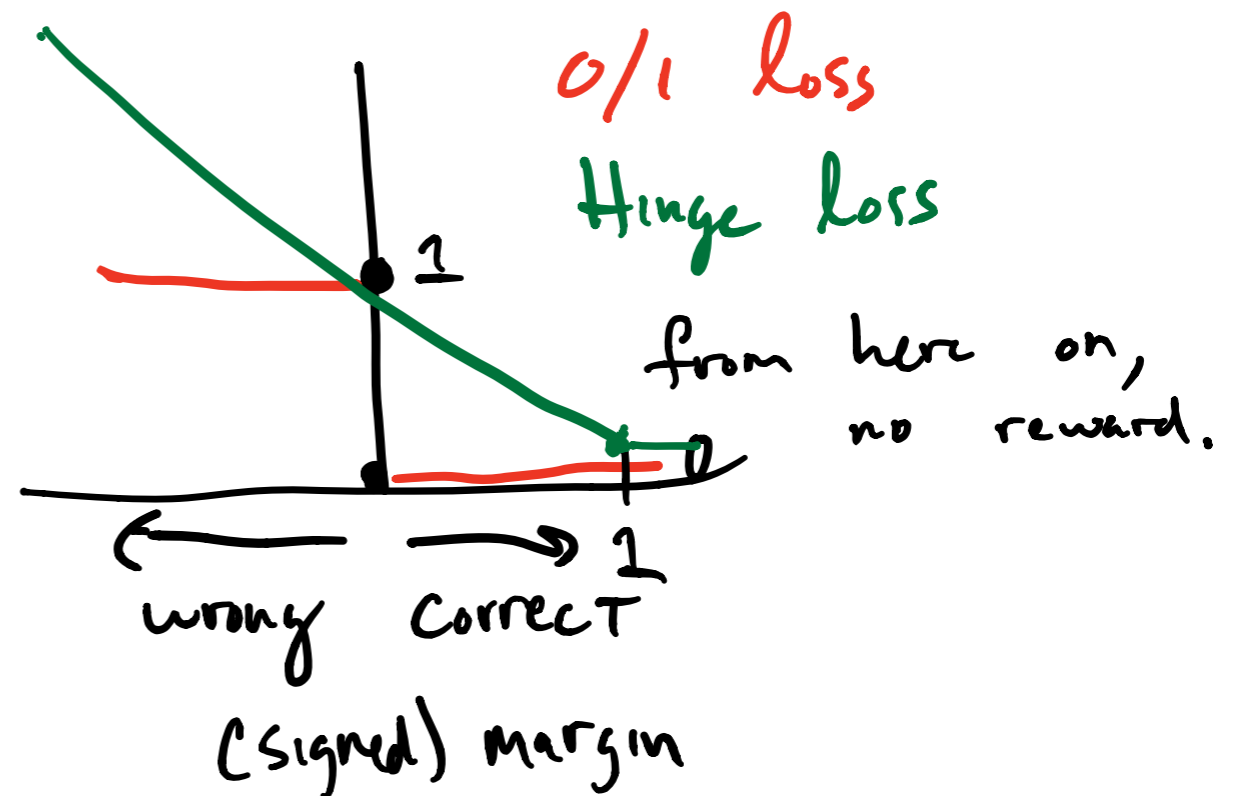
# Smooth loss

Idea: Introduce a "smooth" loss function to make optimization easier

Example: Hinge loss

$$\ell_{Hinge}(y,z) = \text{Max}\{0, 1-y\cdot z\}$$

$y \in \{1, -1\}$

$z = w \cdot x_i$
("raw" output)

0/1 loss

Hinge loss

from here on,
no reward.

1

wrong    correct    1

(signed) margin

# Loss



Zero/one: $\ell^{(0/1)}(y, \hat{y}) = \mathbf{1}[y\hat{y} \leq 0]$

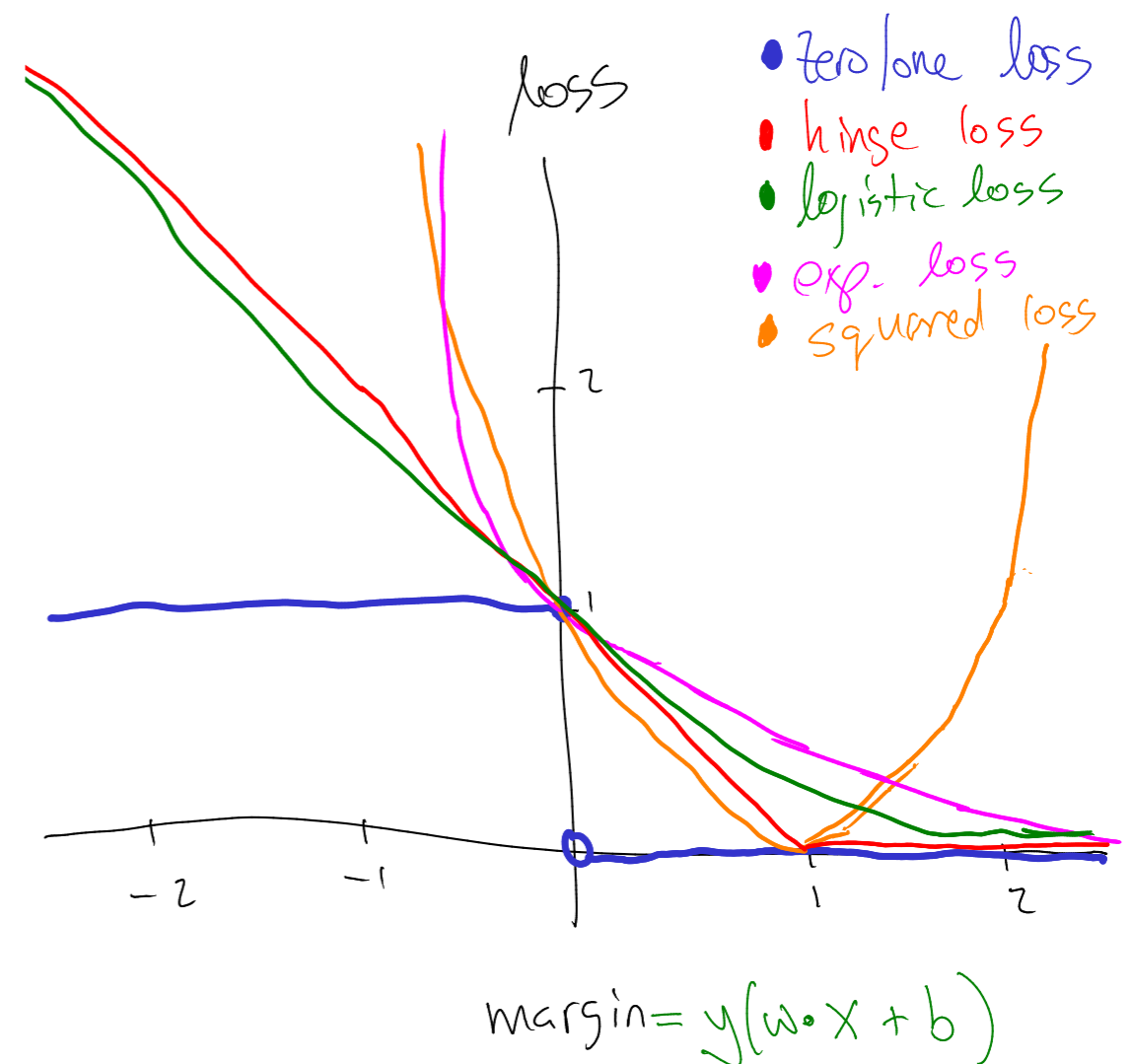Hinge: $\ell^{(\text{hin})}(y, \hat{y}) = \max\{0, 1 - y\hat{y}\}$

Logistic: $\ell^{(\text{log})}(y, \hat{y}) = \frac{1}{\log 2} \log\left(1 + \exp[-y\hat{y}]\right)$

Exponential: $\ell^{(\text{exp})}(y, \hat{y}) = \exp[-y\hat{y}]$

Squared: $\ell^{(\text{sqr})}(y, \hat{y}) = (y - \hat{y})^2$

*Fig and Eq's from CIML [Daume]*

# Regularization

$$\min_{\boldsymbol{w},b} \quad \sum_n \ell(y_n, \boldsymbol{w} \cdot \boldsymbol{x}_n + b)$$

# Regularization

$$\min_{\boldsymbol{w},b} \quad \sum_n \ell(y_n, \boldsymbol{w} \cdot \boldsymbol{x}_n + b) + \boxed{\lambda R(\boldsymbol{w}, b)}$$

Prevent *w* from "getting to crazy"

# Gradient descent

**Algorithm 21** GRADIENTDESCENT($\mathcal{F}, K, \eta_1, \dots$)

1: $z^{(0)} \leftarrow \langle 0, 0, \dots, 0 \rangle$            // initialize variable we are optimizing

2: **for** $k = 1 \dots K$ **do**

3:      $g^{(k)} \leftarrow \nabla_z \mathcal{F}|_{z^{(k\text{-}1)}}$          // compute gradient at current location

4:      $z^{(k)} \leftarrow z^{(k\text{-}1)} - \eta^{(k)} g^{(k)}$          // take a step down the gradient

5: **end for**

6: **return** $z^{(K)}$

$$\nabla_{\boldsymbol{w}} \mathcal{L} = \nabla_{\boldsymbol{w}} \sum_n \exp\left[-y_n(\boldsymbol{w} \cdot \boldsymbol{x}_n + b)\right] + \nabla_{\boldsymbol{w}} \frac{\lambda}{2} \|\boldsymbol{w}\|^2$$

$$\nabla_{\boldsymbol{w}} \mathcal{L} = \nabla_{\boldsymbol{w}} \sum_n \exp\left[-y_n(\boldsymbol{w} \cdot \boldsymbol{x}_n + b)\right] + \nabla_{\boldsymbol{w}} \frac{\lambda}{2} \left\|\boldsymbol{w}\right\|^2$$

$$= \sum_n \left(\nabla_{\boldsymbol{w}} - y_n(\boldsymbol{w} \cdot \boldsymbol{x}_n + b)\right) \exp\left[-y_n(\boldsymbol{w} \cdot \boldsymbol{x}_n + b)\right] + \lambda \boldsymbol{w}$$
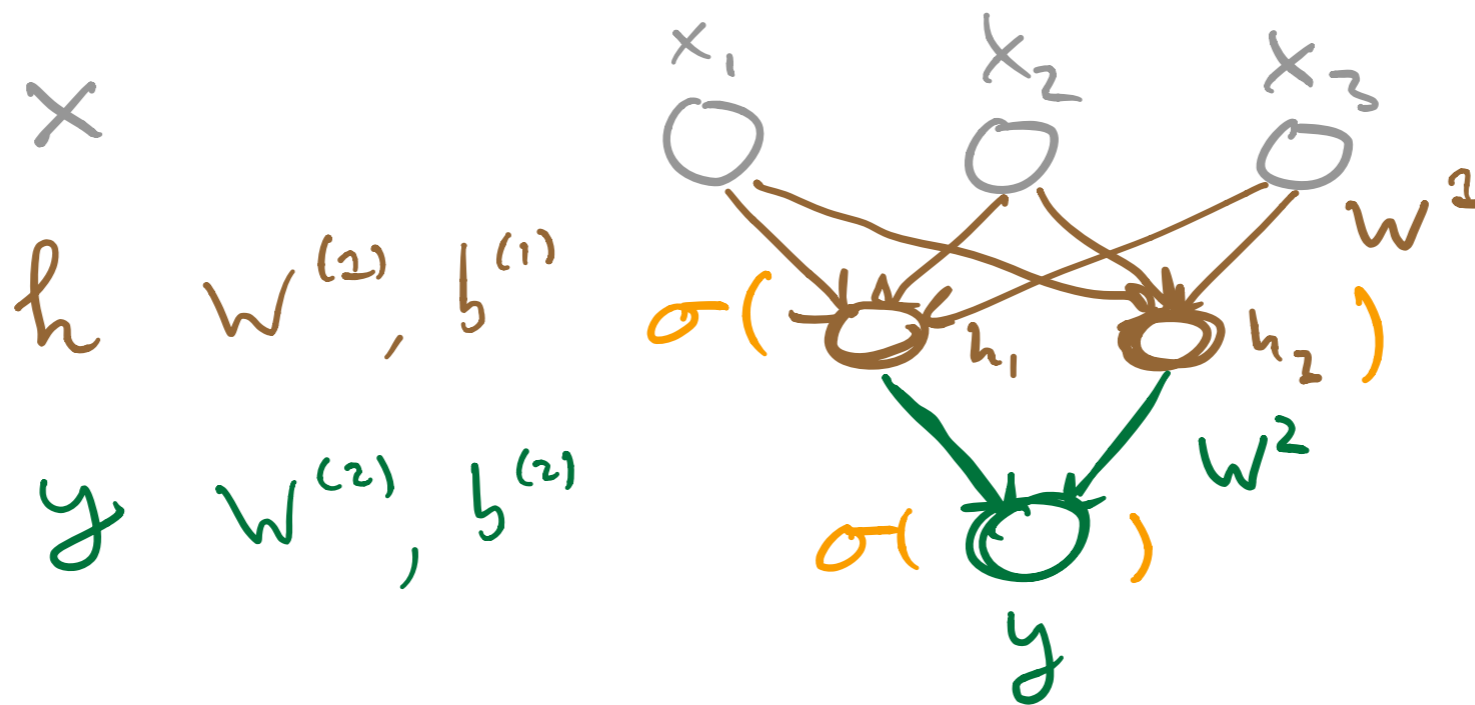
$$\nabla_{\boldsymbol{w}} \mathcal{L} = \nabla_{\boldsymbol{w}} \sum_n \exp\left[-y_n(\boldsymbol{w} \cdot \boldsymbol{x}_n + b)\right] + \nabla_{\boldsymbol{w}} \frac{\lambda}{2} \|\boldsymbol{w}\|^2$$

$$= \sum_n \left(\nabla_{\boldsymbol{w}} - y_n(\boldsymbol{w} \cdot \boldsymbol{x}_n + b)\right) \exp\left[-y_n(\boldsymbol{w} \cdot \boldsymbol{x}_n + b)\right] + \lambda \boldsymbol{w}$$

$$= -\sum_n y_n \boldsymbol{x}_n \exp\left[-y_n(\boldsymbol{w} \cdot \boldsymbol{x}_n + b)\right] + \lambda \boldsymbol{w}$$

# Limitations of linear models

# Neural networks

Idea: Basically stack together a bunch of linear models.

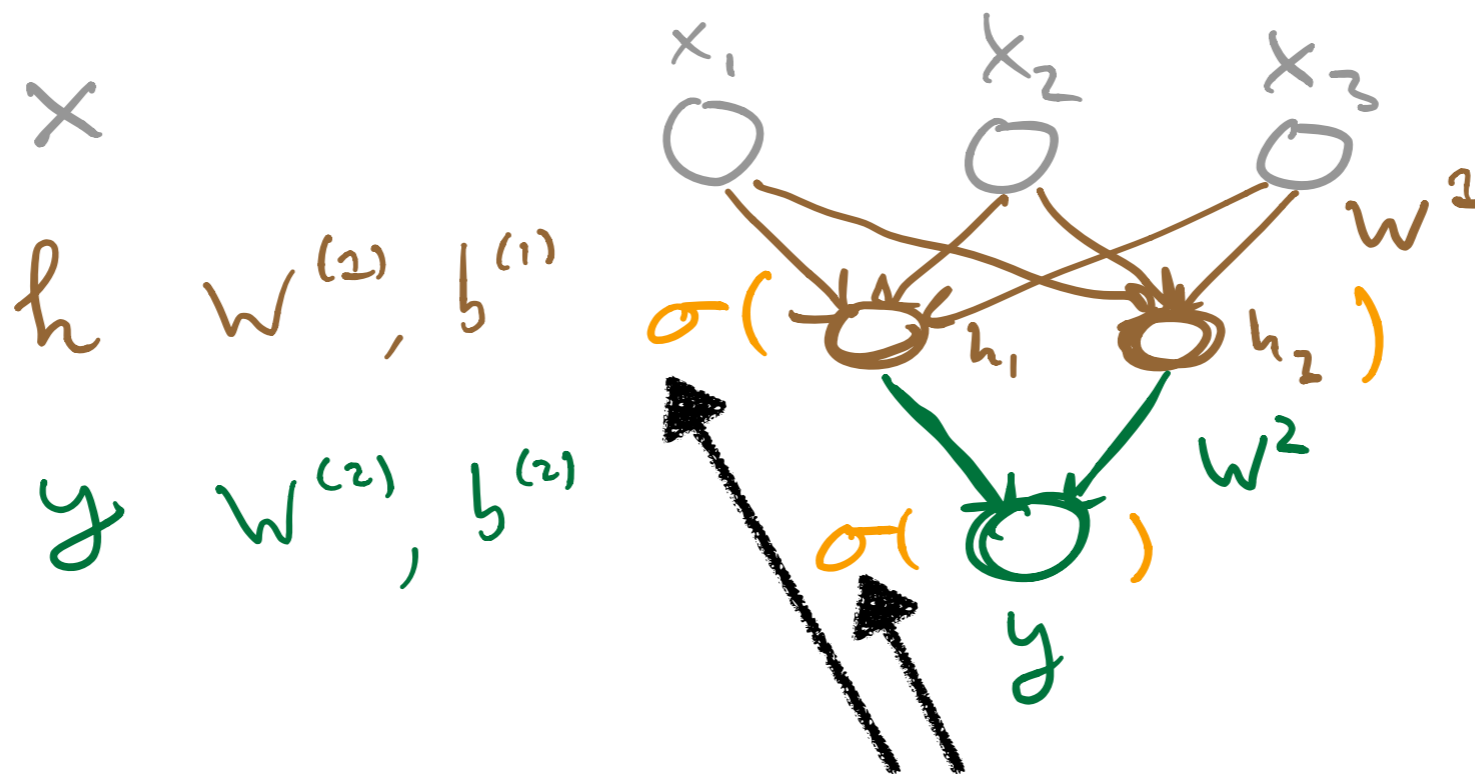This introduces *hidden units* which are neither observations ($x$) nor outputs ($y$)

$x$

$h$ $\quad W^{(1)}, b^{(1)}$

$y$ $\quad W^{(2)}, b^{(2)}$

$x_1$ $\quad x_2$ $\quad x_3$

$\sigma(\quad h_1 \quad h_2)$ $\quad W^1$

$\sigma(\quad) \quad W^2$

$y$

# Neural networks

Idea: Basically stack together a bunch of linear models.

This introduces *hidden units* which are neither observations ($x$) nor outputs ($y$)
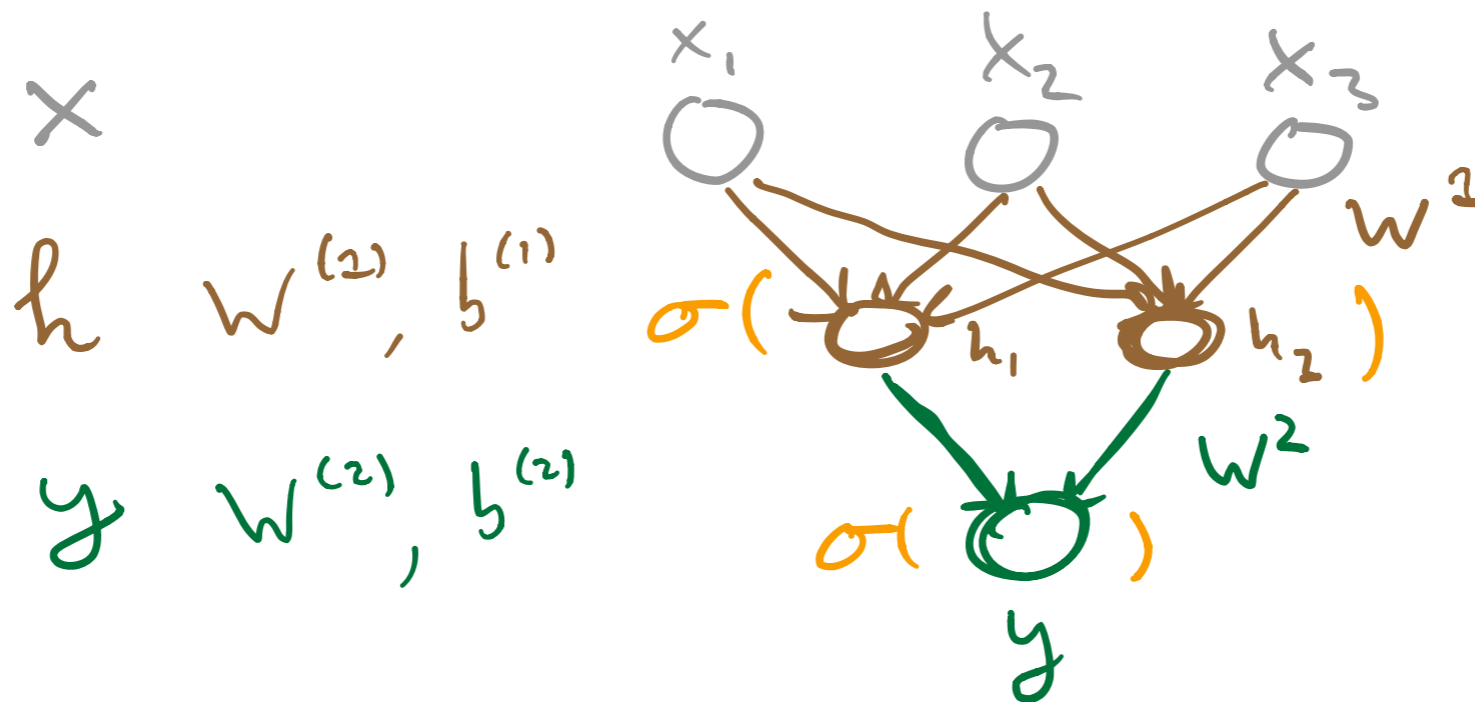


(Non-linear) activation functions

# Neural networks

Idea: Basically stack together a bunch of linear models.

This introduces *hidden units* which are neither observations (*x*) nor outputs (*y*)



The challenge: How do we update weights associated with each node in this *multi-layer* regime?
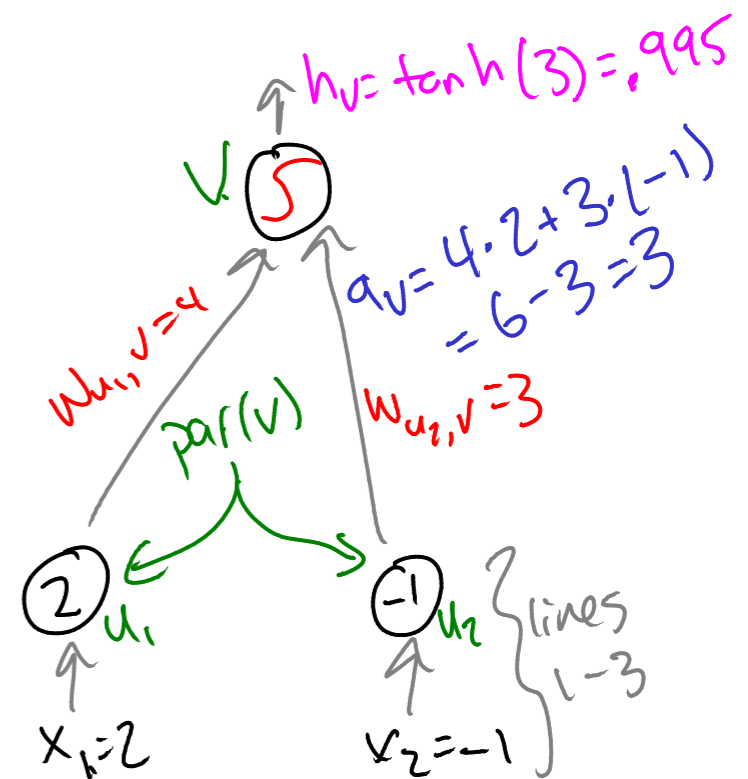
back-propagation = gradient descent + chain rule

---

**Algorithm 27** FORWARDPROPAGATION($x$)

---
1: **for all** input nodes $u$ **do**
2: $\quad h_u \leftarrow$ corresponding feature of $x$
3: **end for**
4: **for all** nodes $v$ in the network whose parent's are computed **do**
5: $\quad a_v \leftarrow \sum_{u \in par(v)} w_{(u,v)} h_u$
6: $\quad h_v \leftarrow \tanh(a_v)$
7: **end for**
8: **return** $a_y$

---

Tanh is another common activation function

$h_v = \tanh(3) = .995$

$V$ ⑤

$a_v = 4 \cdot 2 + 3 \cdot (-1)$
$= 6 - 3 = 3$

$W_{u_1, v} = 4$

$par(v)$

$W_{u_1, v} = 3$

② $u_1$        ⑴ $u_2$  } lines 1-3

$X_1 = 2$        $v_2 = -1$

$y = 2$

$e_y = 2 \cdot 1$
$= 1$        $d_y = 1$   } let $v$ refer to "y"

$W_{u_1, y} = 1$   $par(v)$   $W_{u_2, v} = -1$

$h_y = 3$

$h_{v_2} = 2$

## Handwritten diagram (top right)

$W_{u_1,v}=$ ... $=0$

$par(v)$    $W_{u_2,v}=3$

$(2)\ u_1$     $(-1)\ u_2$   } lines 1-3

$X_1=2$     $V_2=-1$

## Algorithm
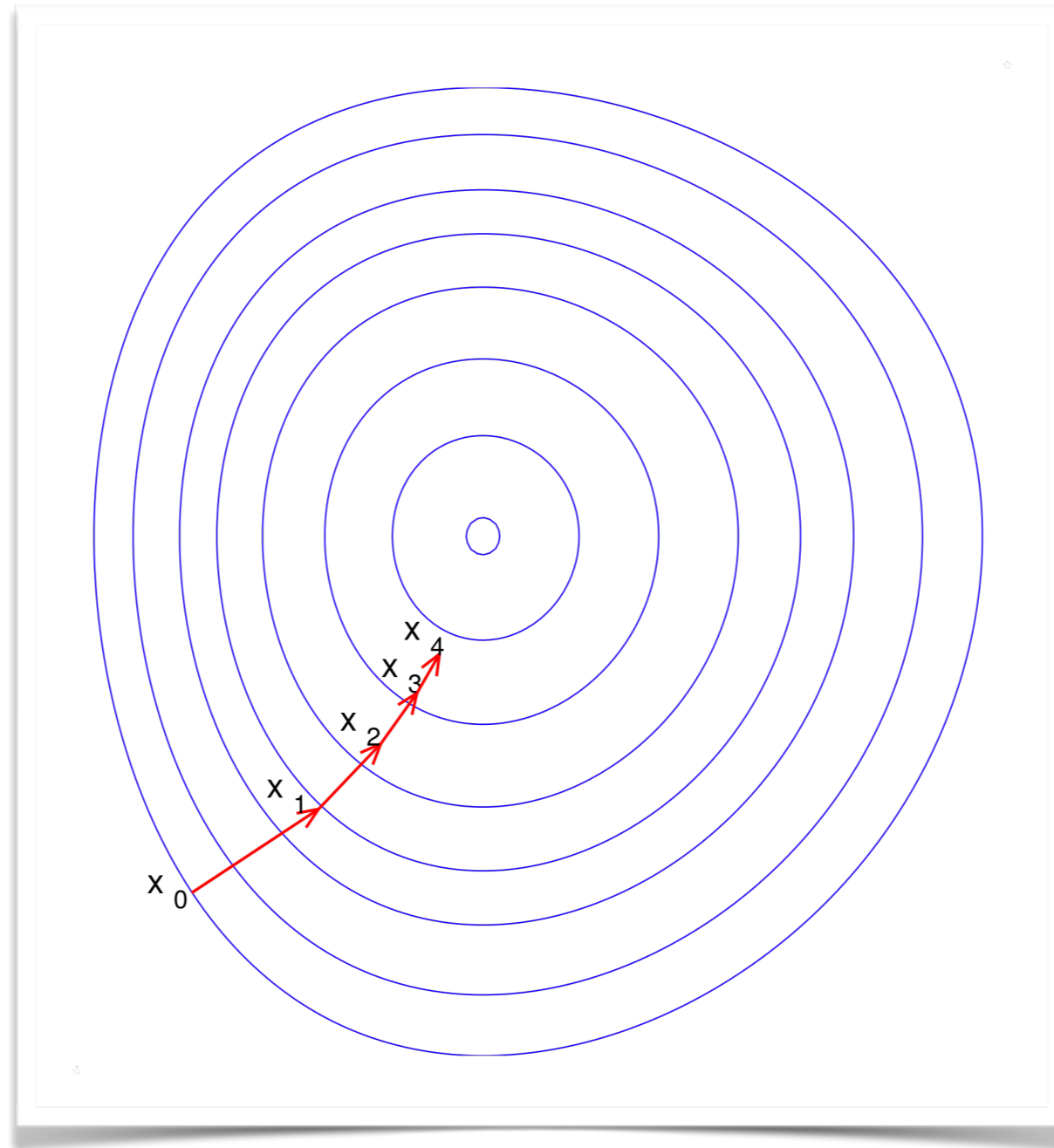
---
**Algorithm 28** BACKPROPAGATION($x$, $y$)

---
1: run FORWARDPROPAGATION($x$) to compute activations
2: $e_y \leftarrow y - a_y$      // compute overall network error
3: **for all** nodes $v$ in the network whose error $e_v$ is computed **do**
4:    **for all** $u \in par(v)$ **do**
5:      $g_{u,v} \leftarrow -e_v h_u$      // compute gradient of this edge
6:      $e_u \leftarrow e_u + e_v w_{u,v}(1 - \tanh^2(a_u))$  // compute the "error" of the parent node
7:    **end for**
8: **end for**
9: **return** all gradients $g_e$

---

## Handwritten diagram (middle right)

$y=2$

$e_y = y$ ... $= 1$    $O\ d_y = 1$   } let $v$ refer to "y"

$h_{v_1}=3\quad W_{u_1,v}=1$    $par(v)$    $W_{u_2,v}=-1\quad h=2$   $v_2$

$e_{v_1}=$ $(1)(1)$ $(1-\tanh^2(a_{u_1}))$

$(v_1)$        $(v_2)\ e_{v_2}=(1)(-1)(1-\tanh^2(a_{u_2}))$

$g_{u,v_1}=-(1)(3)=-3$

$g_{u,v_2}=-(-1)(2)=-2$

## Handwritten diagram (bottom right)

$h_v = \tanh(3) = .995$

$v\ (5)$

$w_{u,v}=4$    $a_v = 4\cdot 2 + 3\cdot(-1)$
     $= 6 - 3 = 3$

What are we doing with these gradients again?

# Gradient descent

# Neural Networks!

If you're interested in learning more…

DS4440 // practical neural networks // spring 2019