

# CSCI567 Machine Learning (Spring 2021)

Sirisha Rambhatla

University of Southern California

Feb 12, 2021

1 / 28

Logistics

## Outline

- 1 Logistics
- 2 Review of last lecture
- 3 Neural Nets

3 / 28

## Outline

- 1 Logistics
- 2 Review of last lecture
- 3 Neural Nets

2 / 28

Logistics

## Logistics

- Sign-up with your group members for the project!

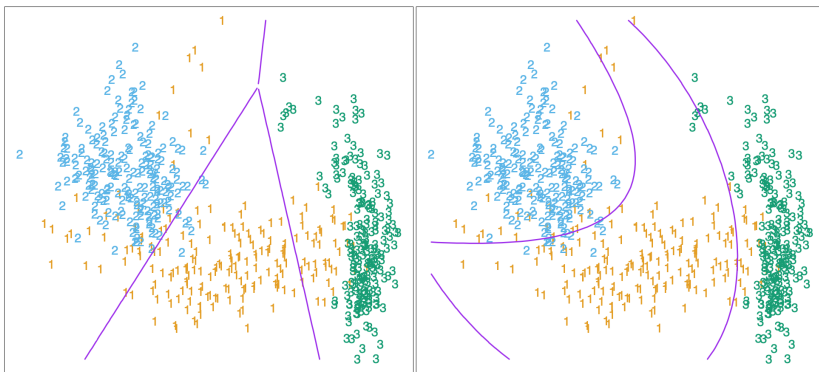
4 / 28

## Outline

- 1 Logistics
- 2 Review of last lecture
- 3 Neural Nets

5 / 28

What do the decision boundaries look like?



The decision boundaries are a quadratic when  $\Sigma$ 's are not the same, this is known as *Quadratic Discriminant Analysis*!

7 / 28

## Linear Discriminant Analysis

The main bottleneck is *not knowing*  $\mathcal{P}(X = \mathbf{x}|y = c)$

$$\mathcal{P}(y = c|X = \mathbf{x}) = \frac{\mathcal{P}(X = \mathbf{x}|y = c)\mathcal{P}(y = c)}{\mathcal{P}(X = \mathbf{x})}.$$

LDA makes two simplifying assumptions:

- Let  $\mathcal{P}(X = \mathbf{x}|y = c) \sim \mathcal{N}(\boldsymbol{\mu}_c, \Sigma_c)$ , and
- Let all class covariances be the same i.e.  $\Sigma_c = \Sigma$  for all  $c \in [C]$

If so, the decision boundary (for binary classification) is given by

$$\begin{aligned} \mathcal{P}(y = 0|X = \mathbf{x}) &= \mathcal{P}(y = 1|X = \mathbf{x}) \\ \mathcal{P}(X = \mathbf{x}|y = 0)\mathcal{P}(y = 0) &= \mathcal{P}(X = \mathbf{x}|y = 1)\mathcal{P}(y = 1) \end{aligned}$$

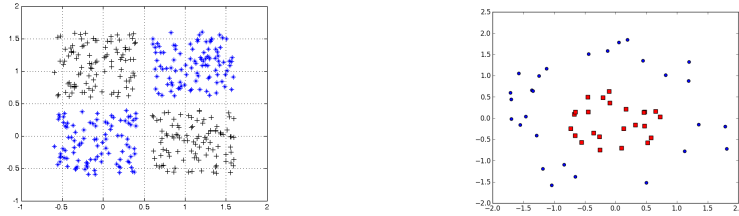
6 / 28

## Outline

- 1 Logistics
- 2 Review of last lecture
- 3 Neural Nets
  - Definition
  - Backpropagation
  - Preventing overfitting

8 / 28

## Linear models are not always adequate



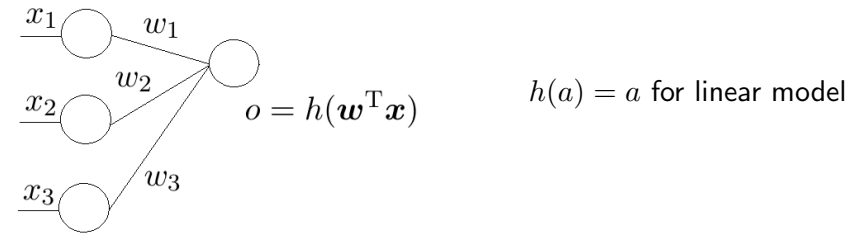
We can use a nonlinear mapping as discussed:

$$\phi(x) : x \in \mathbb{R}^D \rightarrow z \in \mathbb{R}^M$$

*But what kind of nonlinear mapping  $\phi$  should be used? Can we actually learn this nonlinear mapping?*

THE most popular nonlinear models nowadays: **neural nets**

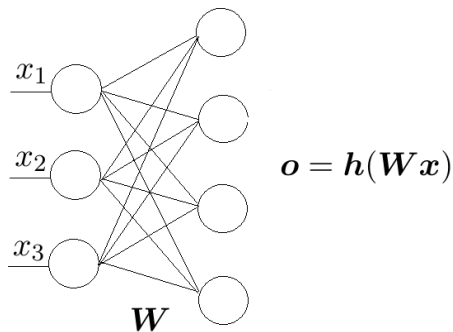
## Linear model as a one-layer neural net



To create non-linearity, can use

- Rectified Linear Unit (**ReLU**):  $h(a) = \max\{0, a\}$
- sigmoid function:  $h(a) = \frac{1}{1+e^{-a}}$
- TanH:  $h(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$
- many more

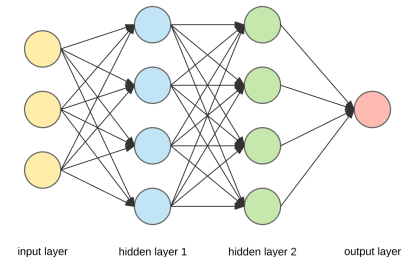
## More output nodes



$$W \in \mathbb{R}^{4 \times 3}, h : \mathbb{R}^4 \rightarrow \mathbb{R}^4 \text{ so } h(a) = (h_1(a_1), h_2(a_2), h_3(a_3), h_4(a_4))$$

Can think of this as a nonlinear basis:  $\Phi(x) = h(Wx)$

## More layers



Becomes a network:

- each node is called a **neuron**
- $h$  is called the **activation function**
  - can use  $h(a) = 1$  for one neuron in each layer to *incorporate bias term*
  - output neuron can use  $h(a) = a$
- #layers refers to #hidden\_layers (plus 1 or 2 for input/output layers)
- **deep** neural nets can have many layers and *millions* of parameters
- this is a **feedforward, fully connected** neural net, there are many variants

## How powerful are neural nets?

**Universal approximation theorem** (Cybenko, 89; Hornik, 91):

*A feedforward neural net with a single hidden layer can approximate any continuous functions.*

It might need a huge number of neurons though, and *depth helps!*

Designing network architecture is important and very complicated

- for feedforward network, need to decide number of hidden layers, number of neurons at each layer, activation functions, etc.

13 / 28

## Learning the model

*No matter how complicated the model is, our goal is the same:* minimize

$$\mathcal{E}(\mathbf{W}_1, \dots, \mathbf{W}_L) = \frac{1}{N} \sum_{n=1}^N \mathcal{E}_n(\mathbf{W}_1, \dots, \mathbf{W}_L)$$

where

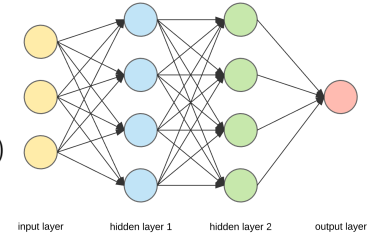
$$\mathcal{E}_n(\mathbf{W}_1, \dots, \mathbf{W}_L) = \begin{cases} \|f(\mathbf{x}_n) - \mathbf{y}_n\|_2^2 & \text{for regression} \\ \ln \left( 1 + \sum_{k \neq y_n} e^{f(\mathbf{x}_n)_k - f(\mathbf{x}_n)_{y_n}} \right) & \text{for classification} \end{cases}$$

15 / 28

## Math formulation

An L-layer neural net can be written as

$$f(\mathbf{x}) = \mathbf{h}_L(\mathbf{W}_L \mathbf{h}_{L-1}(\mathbf{W}_{L-1} \cdots \mathbf{h}_1(\mathbf{W}_1 \mathbf{x})))$$



To ease notation, for a given input  $\mathbf{x}$ , define recursively

$$\mathbf{o}_0 = \mathbf{x}, \quad \mathbf{a}_\ell = \mathbf{W}_\ell \mathbf{o}_{\ell-1}, \quad \mathbf{o}_\ell = \mathbf{h}_\ell(\mathbf{a}_\ell) \quad (\ell = 1, \dots, L)$$

where

- $\mathbf{W}_\ell \in \mathbb{R}^{D_\ell \times D_{\ell-1}}$  is the weights between layer  $\ell - 1$  and  $\ell$
- $D_0 = D, D_1, \dots, D_L$  are numbers of neurons at each layer
- $\mathbf{a}_\ell \in \mathbb{R}^{D_\ell}$  is input to layer  $\ell$
- $\mathbf{o}_\ell \in \mathbb{R}^{D_\ell}$  is output to layer  $\ell$
- $\mathbf{h}_\ell : \mathbb{R}^{D_\ell} \rightarrow \mathbb{R}^{D_\ell}$  is activation functions at layer  $\ell$

14 / 28

## How to optimize such a complicated function?

Same thing: apply **SGD**! even if the model is *nonconvex*.

What is the gradient of this complicated function?

*Chain rule is the only secret:*

- for a composite function  $f(g(w))$

$$\frac{\partial f}{\partial w} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial w}$$

- for a composite function  $f(g_1(w), \dots, g_d(w))$

$$\frac{\partial f}{\partial w} = \sum_{i=1}^d \frac{\partial f}{\partial g_i} \frac{\partial g_i}{\partial w}$$

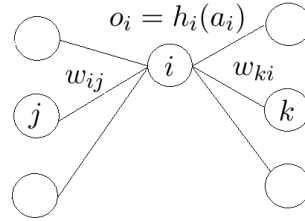
the simplest example  $f(g_1(w), g_2(w)) = g_1(w)g_2(w)$

16 / 28

## Computing the derivative

Drop the subscript  $\ell$  for layer for simplicity.

Find the **derivative of  $\mathcal{E}_n$  w.r.t. to  $w_{ij}$**



$$\frac{\partial \mathcal{E}_n}{\partial w_{ij}} = \frac{\partial \mathcal{E}_n}{\partial a_i} \frac{\partial a_i}{\partial w_{ij}} = \frac{\partial \mathcal{E}_n}{\partial a_i} \frac{\partial (w_{ij} o_j)}{\partial w_{ij}} = \frac{\partial \mathcal{E}_n}{\partial a_i} o_j$$

$$\frac{\partial \mathcal{E}_n}{\partial a_i} = \frac{\partial \mathcal{E}_n}{\partial o_i} \frac{\partial o_i}{\partial a_i} = \left( \sum_k \frac{\partial \mathcal{E}_n}{\partial a_k} \frac{\partial a_k}{\partial o_i} \right) h'_i(a_i) = \left( \sum_k \frac{\partial \mathcal{E}_n}{\partial a_k} w_{ki} \right) h'_i(a_i)$$

17 / 28

## Computing the derivative

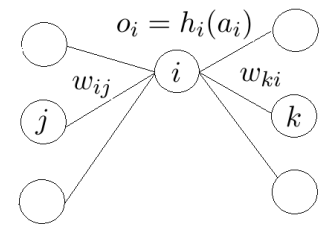
Adding the subscript for layer:

$$\frac{\partial \mathcal{E}_n}{\partial w_{\ell,ij}} = \frac{\partial \mathcal{E}_n}{\partial a_{\ell,i}} o_{\ell-1,j}$$

$$\frac{\partial \mathcal{E}_n}{\partial a_{\ell,i}} = \left( \sum_k \frac{\partial \mathcal{E}_n}{\partial a_{\ell+1,k}} w_{\ell+1,ki} \right) h'_{\ell,i}(a_{\ell,i})$$

For the last layer, for square loss

$$\frac{\partial \mathcal{E}_n}{\partial a_{L,i}} = \frac{\partial (h_{L,i}(a_{L,i}) - y_{n,i})^2}{\partial a_{L,i}} = 2(h_{L,i}(a_{L,i}) - y_{n,i}) h'_{L,i}(a_{L,i})$$



18 / 28

## Computing the derivative

Using **matrix notation** greatly simplifies presentation and implementation:

$$\frac{\partial \mathcal{E}_n}{\partial \mathbf{W}_\ell} = \frac{\partial \mathcal{E}_n}{\partial \mathbf{a}_\ell} \mathbf{o}_{\ell-1}^T$$

$$\frac{\partial \mathcal{E}_n}{\partial \mathbf{a}_\ell} = \begin{cases} \left( \mathbf{W}_{\ell+1}^T \frac{\partial \mathcal{E}_n}{\partial \mathbf{a}_{\ell+1}} \right) \circ \mathbf{h}'_\ell(\mathbf{a}_\ell) & \text{if } \ell < L \\ 2(\mathbf{h}_L(\mathbf{a}_L) - \mathbf{y}_n) \circ \mathbf{h}'_L(\mathbf{a}_L) & \text{else} \end{cases}$$

where  $\mathbf{v}_1 \circ \mathbf{v}_2 = (v_{11}v_{21}, \dots, v_{1D}v_{2D})$  is the element-wise product (a.k.a. Hadamard product).

Verify yourself!

19 / 28

## Putting everything into SGD

The **backpropagation** algorithm (**Backprop**)

Initialize  $\mathbf{W}_1, \dots, \mathbf{W}_L$ . Repeat:

- ① randomly pick one data point  $n \in [N]$
- ② **forward propagation**: for each layer  $\ell = 1, \dots, L$ 
  - compute  $\mathbf{a}_\ell = \mathbf{W}_\ell \mathbf{o}_{\ell-1}$  and  $\mathbf{o}_\ell = \mathbf{h}_\ell(\mathbf{a}_\ell)$  ( $\mathbf{o}_0 = \mathbf{x}_n$ )
- ③ **backward propagation**: for each  $\ell = L, \dots, 1$ 
  - compute

$$\frac{\partial \mathcal{E}_n}{\partial \mathbf{a}_\ell} = \begin{cases} \left( \mathbf{W}_{\ell+1}^T \frac{\partial \mathcal{E}_n}{\partial \mathbf{a}_{\ell+1}} \right) \circ \mathbf{h}'_\ell(\mathbf{a}_\ell) & \text{if } \ell < L \\ 2(\mathbf{h}_L(\mathbf{a}_L) - \mathbf{y}_n) \circ \mathbf{h}'_L(\mathbf{a}_L) & \text{else} \end{cases}$$

- update weights

$$\mathbf{W}_\ell \leftarrow \mathbf{W}_\ell - \eta \frac{\partial \mathcal{E}_n}{\partial \mathbf{W}_\ell} = \mathbf{W}_\ell - \eta \frac{\partial \mathcal{E}_n}{\partial \mathbf{a}_\ell} \mathbf{o}_{\ell-1}^T$$

*Think about how to do the last two steps properly!*

20 / 28

## More tricks to optimize neural nets

Many variants based on backprop

- SGD with **minibatch**: randomly sample a batch of examples to form a stochastic gradient
- SGD with **momentum**
- ...

21 / 28

## SGD with momentum

Initialize  $w_0$  and **velocity**  $v = 0$

For  $t = 1, 2, \dots$

- form a stochastic gradient  $g_t$
- update velocity  $v \leftarrow \alpha v - \eta g_t$  for some discount factor  $\alpha \in (0, 1)$
- update weight  $w_t \leftarrow w_{t-1} + v$

Updates for first few rounds:

- $w_1 = w_0 - \eta g_1$
- $w_2 = w_1 - \alpha \eta g_1 - \eta g_2$
- $w_3 = w_2 - \alpha^2 \eta g_1 - \alpha \eta g_2 - \eta g_3$
- ...

22 / 28

## Overfitting

**Overfitting is very likely** since the models are too powerful.

Methods to overcome overfitting:

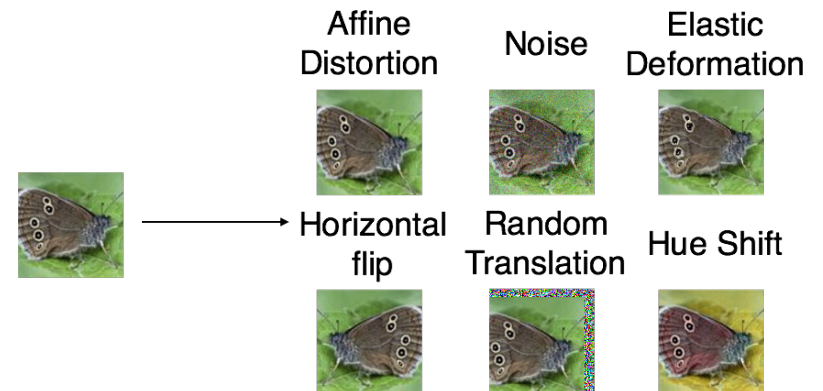
- data augmentation
- regularization
- dropout
- early stopping
- ...

23 / 28

## Data augmentation

Data: the more the better. How do we get more data?

**Exploit prior knowledge to add more training data**



24 / 28

## Regularization

**L2 regularization:** minimize

$$\mathcal{E}'(\mathbf{W}_1, \dots, \mathbf{W}_L) = \mathcal{E}(\mathbf{W}_1, \dots, \mathbf{W}_L) + \lambda \sum_{\ell=1}^L \|\mathbf{W}_\ell\|_2^2$$

Simple change to the gradient:

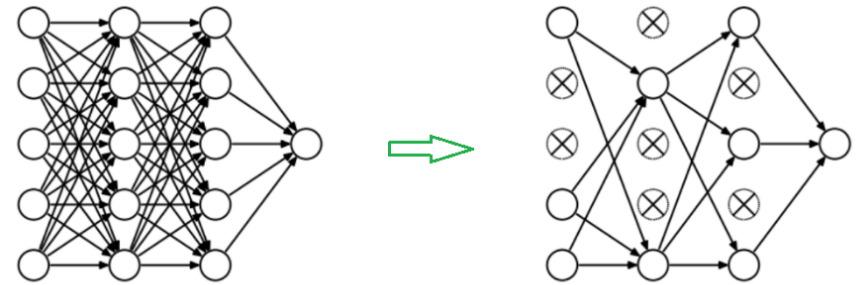
$$\frac{\partial \mathcal{E}'}{\partial w_{ij}} = \frac{\partial \mathcal{E}}{\partial w_{ij}} + 2\lambda w_{ij}$$

Introduce *weight decaying effect*

25 / 28

## Dropout

**Randomly delete neurons** during training

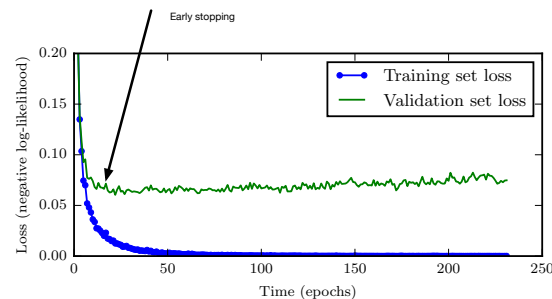


Very effective, makes training faster as well

26 / 28

## Early stopping

Stop training when the performance on validation set stops improving



27 / 28

## Conclusions for neural nets

Deep neural networks

- are hugely popular, achieving *best performance* on many problems
- do need *a lot of data* to work well
- take *a lot of time* to train (need GPUs for massive parallel computing)
- take some work to select architecture and hyperparameters
- are still not well understood in theory

28 / 28