

Lecture 14

following Higham + Higham, SIAM Rev 2020

Deep Learning + the Stochastic Gradient Method

Sigmoid function:  $\hat{\sigma}(x) = \frac{1}{1+e^{-x}}$

Fig 2.2.  $\hat{\sigma}: \mathbb{R} \rightarrow \mathbb{R}$ 

smoothed version of a step function.

$$\hat{\sigma}'(x) = \frac{-1}{(1+e^{-x})^2} (-e^{-x}) =$$

$$= \hat{\sigma}(x) \left( \frac{+e^{-x}}{1+e^{-x}} \right) = \hat{\sigma}(x) \left( \frac{1+e^{-x}-1}{1+e^{-x}} \right)$$

$$= \hat{\sigma}(x) (1 - \hat{\sigma}(x))$$

$\hat{\sigma}(x)$  is a smoothed version of the STEP function which mimics a single neuron in the brain: firing, or not firing.

scale + shift (weight + bias in NN language)

$$\hat{\sigma}(w(x+b))$$

Vector version:  $\sigma: \mathbb{R}^m \rightarrow \mathbb{R}^m$ 

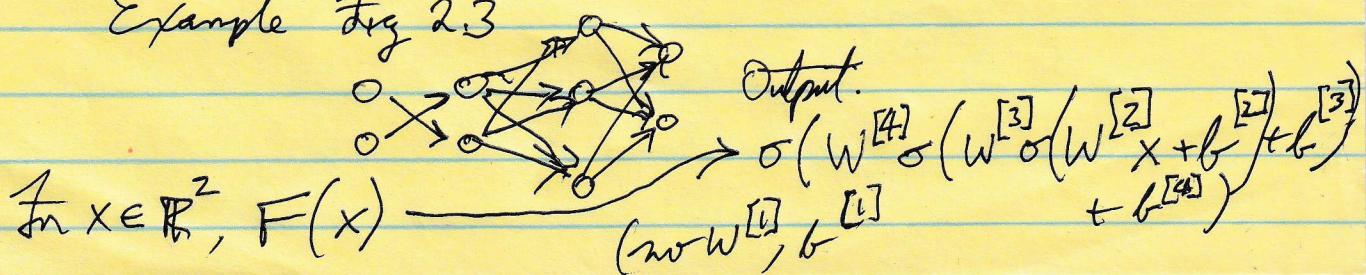
$$(\sigma(z))_i = \hat{\sigma}(z_i) \quad z \in \mathbb{R}^m$$

We want Layers of neurons:output from one layer =  $a$  = input to next layeroutput from next layer =  $\sigma(Wa+b)$ 

$$W \in \mathbb{R}^{m \times m} \quad b \in \mathbb{R}^m$$

square iff <sup>t</sup># neurons in input layer  
- # " " output for layer.

Example Fig 2.3

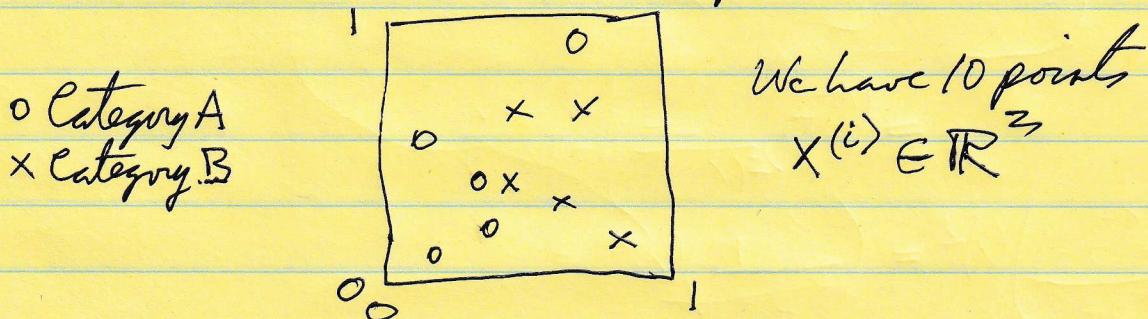


DL2

$$\begin{array}{lll} W^{[2]} \in \mathbb{R}^{2 \times 2} & b^{[2]} \in \mathbb{R}^2 & \# \text{pars} = 6 \\ W^{[3]} \in \mathbb{R}^{3 \times 2} & b^{[3]} \in \mathbb{R}^3 & \# \text{pars} = 9 \\ W^{[4]} \in \mathbb{R}^{2 \times 3} & b^{[4]} \in \mathbb{R}^2 & \# \text{pars} = 8 \end{array}$$

total 23

We'd like to build a "classifier" based on data



We'll arbitrarily say we'd like the output  $F(x)$  to be close to  $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$  for points  $x$  in cat. A and  $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$  for points  $x$  in cat. B. (These vectors  $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$  and  $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$  have no relationship to the  $x$ -picture.) We write

$$y(x^{(i)}) = \begin{cases} \begin{bmatrix} 1 \\ 0 \end{bmatrix} & \text{if } x^{(i)} \text{ is cat A} \\ \begin{bmatrix} 0 \\ 1 \end{bmatrix} & \text{if } x^{(i)} \text{ is cat B.} \end{cases}$$

We want to choose the 23 parameters to

$$\begin{aligned} &\text{minimize Cost}(W^{[2]}, W^{[3]}, W^{[4]}, b^{[2]}, b^{[3]}, b^{[4]}) \\ &= \frac{1}{10} \sum_{i=1}^{10} \frac{1}{2} \| y(x^{(i)}) - F(x^{(i)}) \|_2^2 \end{aligned}$$

defined on p. DL1

This is a nonlinear, nonconvex opt problem — called a nonlinear least squares problem.

Possible Methods to Minimize Cost

BFGS

Newton

Barzilai-Borwein

Levenberg-Maquardt } implemented in lsqnonlin  
Trust region method } (Matlab opt toolbox).

DL 3

Using algorithm to choose the 23 parameters —  
TRAINING THE NETWORK, + then plotting the points in  
the  $\times$  square where  $F_1(x) > F_2(x)$  (meaning: best  
choice for  $x$  is cat. A) gives Fig 2.4.  
Adding one more cat. B point, + re-optimizing  
(re-training), get Fig. 2.5.

However, there is no guarantee any method will  
find the globally optimal choice of parameters —  
especially if there are a lot of them.

Deep Learning  $\equiv$  Many layers in NN.

So far, we are not exploiting the property that the  
output from one layer is a simple function of the input  
to it.

Assume there are  $L$  layers labelled  $l=1, \dots, l=L$ ,  
with  $n_l$  neurons in layer  $l$ .  
Let  $a_j^{(l)}$  = output, or activation, from neuron  $j$  at  
layer  $l$ . We have

$$a^{[1]} = x = \text{input}$$
$$a^{[l]} = \sigma(W^{[l]} a^{[l-1]} + b^{[l]}) \in \mathbb{R}^{n_l}, \quad l=2, 3, \dots, L$$

"feed forward". Now want to minimize  
 $\text{Cost}(W^{[2]}, \dots, W^{[L]}, b^{[2]}, \dots, b^{[L]})$

$$C = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} \|y(x^{(i)}) - a^{[L]}(x^{(i)})\|_2^2$$

Write  $p^T = ((\text{vec}(W^{[2]}))^T, \dots, (\text{vec}(W^{[L]}))^T, (b^{[2]})^T, \dots, (b^{[L]})^T)$   
vector of parameters  $p \in \mathbb{R}^s$

## Stochastic Gradient Method. (SG "D")

Two key differences from the gradient method we used before.

(1) Replace  $\nabla \text{Cost}(\rho)$  by  $\nabla (\text{only } \cancel{\text{one or a few}} \text{ terms in Cost}(\rho))$

These are chosen randomly

(2) No line search. Update  $\rho \leftarrow \rho + \eta \nabla (\cdot)$ ,

where  $\eta$  is small (stepsize or "learning rate" in ML)

$$\text{let } C_{x^{(i)}} = \frac{1}{2} \| y(x^{(i)}) - \underbrace{a^{[L]}(x^{(i)})}_{\text{depends on } \rho} \|_2^2 \equiv C_{x^{(i)}}(\rho)$$

Then  $C = \frac{1}{N} \sum_{i=1}^N C_{x^{(i)}}$  and depends on  $\rho$ .

$$\nabla C(\rho) = \frac{1}{N} \sum_{i=1}^N \nabla C_{x^{(i)}}(\rho)$$

### Stochastic Gradient Method

Initialize  $\rho$ . For  $k=1, 2, \dots$

Choose  $i$  randomly from  $\{1, 2, 3, \dots, N\}$   
Update  $\rho \leftarrow \rho - \eta \nabla C_{x^{[i]}}(\rho)$

called "sampling with replacement"  
because the same index might  
be sampled multiple times.

Sampling without replacement: pick a permutation of  $\{1, 2, \dots, N\}$  randomly & then run through all indices once.  
This is called an epoch, & may be repeated multiple times.

More common: instead of one index at a time, choose  $m$  integers  $k_1, \dots, k_m$  randomly from  $\{1, \dots, N\}$  and use

$$\rho \leftarrow \rho - \eta \frac{1}{m} \sum_{i=1}^m \nabla C_{x^{[k_i]}}(\rho).$$

This is a MINIBATCH.

DL5

In SG  
For FIXED  $C$

## Gradients via Back Propagation (Automatic Differentiation)

→ We need  $\frac{\partial C}{\partial w_{jk}^{[l]}}$  and  $\frac{\partial C}{\partial b_j^{[l]}}$  as these make up  $\nabla C(p)$ .

Key point: the training data  $x^{(i)}$  and  $y(x^{(i)})$  are fixed.

(A) So write  $C = \frac{1}{2} \|y - a^{[L]}\|_2^2$   $a^{[L]}(x^{(i)})$

$C, x^{(i)}, y(x^{(i)})$  depends on  $w_{jk}^{[l]}, b_j^{[l]}$  (weights & biases)

(B) Let  $z^{[l]} = w^{[l]}a^{[l-1]} + b^{[l]} \in \mathbb{R}^m$   $l=2, 3, \dots, L$

$z_j^{[l]}$  is weighted input for neuron  $j$  at layer  $l$ .

(C) Then  $a^{[l]} = \sigma(z^{[l]})$

(D) Let  $\delta_j^{(l)} = \frac{\partial C}{\partial z_j^{(l)}}$  NOT composition as in Sec 12.

Notation Let  $p, q \in \mathbb{R}^n$ , then  $p \circ q \in \mathbb{R}^n$  is defined by  $(p \circ q)_i = p_i q_i$  componentwise product.

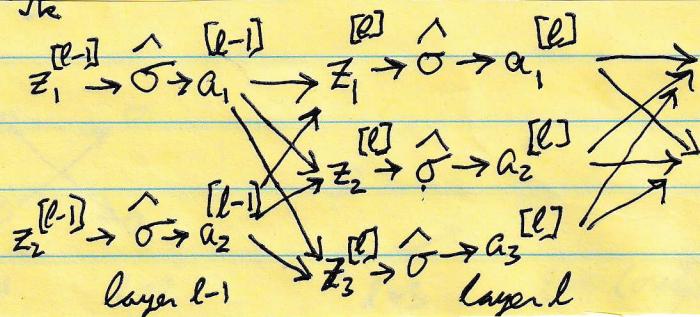
Recall  
i fixed

Lemma

$\frac{\partial C}{\partial z_j^{[l]}}$  (i)  $\delta^{[L]} = \sigma'(z^{[L]}) \circ (a^{[L]} - y)$   
 $\frac{\partial C}{\partial z_j^{[l]}}$  (ii)  $\delta^{[l]} = \sigma'(z^{[l]}) \circ ((W^{[l+1]})^T \delta^{[l+1]})$   $2 \leq l \leq L-1$   
 $j=1, 2, \dots$   
 $\frac{\partial C}{\partial b_j^{[l]}}$  (iii)  $\delta_j^{[l]} = \delta_j^{[l]}$   $2 \leq l \leq L$

WHAT  
WE  
NEED  
FOR  $\nabla C(p)$

(iv)  $\frac{\partial C}{\partial w_{jk}^{[l]}} = \delta_j^{[l]} a_k^{[l-1]}$   $2 \leq l \leq L$



$$z_l = W a^{[l-1]} + b^{[l]}$$

DL6

Recall  
i is fixed

Proof

(i) We have  $a_{(c)}^{[l]} = \sigma(z^{[l]})$  so  $\frac{\partial a_{ij}^{[l]}}{\partial z_j^{[l]}} = \hat{\sigma}'(z_j^{[l]})$

$$\text{and } (A) \quad \frac{\partial C}{\partial a_{ij}^{[l]}} = \frac{\partial}{\partial a_{ij}^{[l]}} \frac{1}{2} \sum_{k=1}^{n_L} (y_k - a_{ik}^{[l]})^2 = -(y_j - a_{ij}^{[l]})$$

so by the ordinary chain rule

$$s_j^{[l]} = \frac{\partial C}{\partial z_j^{[l]}} = \frac{\partial C}{\partial a_{ij}^{[l]}} \frac{\partial a_{ij}^{[l]}}{\partial z_j^{[l]}} = (a_{ij}^{[l]} - y_j) \hat{\sigma}'(z_j^{[l]})$$

which is equivalent to (i).

again by  
chain rule,  
for  $l < L$ ,

$$(ii) \quad s_j^{[l]} = \frac{\partial C}{\partial z_j^{[l]}} = \sum_{k=1}^{n_{l+1}} \frac{\partial C}{\partial z_k^{[l+1]}} \left( \frac{\partial z_k^{[l+1]}}{\partial z_j^{[l]}} \right) = \sum_{k=1}^{n_{l+1}} s_k^{[l+1]} \frac{\partial z_k^{[l+1]}}{\partial z_j^{[l]}}$$

The  $z_k$  at layer  $l+1$  are connected to  $z_j$  at layer  $l$  by

$$z_k^{[l+1]} = \sum_{s=1}^{n_l} w_{ks} \sigma(z_s^{[l]}) + b_k^{[l+1]}$$

Need to consider how all  $z_k^{[l+1]}$  depend on the particular  $z_j^{[l]}$  that we are differentiating  $C$  with respect to.

$$\text{so } \frac{\partial z_k^{[l+1]}}{\partial z_j^{[l]}} = w_{kj} \hat{\sigma}'(z_j^{[l]})$$

$$\text{so becomes } s_j^{[l]} = \hat{\sigma}'(z_j^{[l]}) \left( (W^{[l+1]})^T s^{[l+1]} \right) = (ii)$$

(iii) We have  $z_j^{[l]} = (B) (W^{[l]} a^{[l-1]} + b^{[l]}) = (B) (W^{[l]} \sigma(z^{[l-1]})) + b_j^{[l]}$

since  $z^{[l-1]}$  does not depend on  $b_j^{[l]}$  this says

$$\frac{\partial z_j^{[l]}}{\partial b_j^{[l]}} = 1$$

$$\text{so } \frac{\partial C}{\partial \theta_j^{[l]}} = \frac{\partial C}{\partial z_j^{[l]}} \frac{\partial z_j^{[l]}}{\partial \theta_j^{[l]}} \stackrel{(i)}{=} \delta_j^{[l]} \quad \checkmark$$

$$(iv) \text{ From (B), } z_j^{[l]} = \sum_{k=1}^{m_{l-1}} w_{jk}^{[l]} a_k^{[l-1]} + b_j^{[l]}$$

$$\text{so } \frac{\partial z_j^{[l]}}{\partial w_{jk}^{[l]}} = a_k^{[l-1]} \quad \forall j \quad (\text{jth neuron at layer } l \text{ uses weights only in } j^{\text{th}} \text{ row of } W^{[l]})$$

and

$$\frac{\partial z_s^{[l]}}{\partial w_{jk}^{[l]}} = 0 \quad \text{for } s \neq j.$$

$$\text{so } \frac{\partial C}{\partial w_{jk}^{[l]}} = \sum_{s=1}^{m_l} \frac{\partial C}{\partial z_s^{[l]}} \frac{\partial z_s^{[l]}}{\partial w_{jk}^{[l]}} \stackrel{(ii)}{=} \left( \frac{\partial C}{\partial z_j^{[l]}} a_k^{[l-1]} + 0 \right)$$

$\delta_j^{[l]} \xrightarrow{(ii)}$

### Key Points of the Lemma:

- $a^{[l]}$  is computed by forward pass through the network, computing  $a^{[1]}, a^{[2]}, \dots, a^{[L]}$  one at a time.
  - the (i) gives us  $\delta^{[L]}$
  - the (ii) gives us  $\delta^{[L-1]}, \delta^{[L-2]}, \dots, \delta^{[2]}$  in a backward pass through the network. This is called backpropagation. It is a special case of the REVERSE MODE of AUTOMATIC DIFFERENTIATION (see N+W Ch. 8.2) (See also forward+backward).
- (idea goes back many decades)

Can replace Hadamard notation by diagonal matrices.

$$\text{Let } D^{[l]} = \text{Diag}(\delta'(z_1^{[l]}), \dots, \delta'(z_{m_l}^{[l]})). \text{ Then } \stackrel{(ii)}{\delta^{[l]}} = D^{[l]} (W^{[l+1]})^T \delta^{[l+1]} \quad 2 \leq l \leq L-1$$

DLP

Pseudocode : SG with BP

hence  $\leq G$

for counter = 1, 2, ...

choose k randomly from  $\{1, 2, \dots, N\}$  (a.k.a. i)

$a^{[1]} = x^{(k)}$   $\rightarrow$  for  $l = 2$  to L

$$z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$$

$$a^{[l]} = \sigma(z^{[l]})$$

$$D^{[l]} = \text{Diag}(\sigma'(z^{[l]}))$$

end

$$\delta^{[L]} = D^{[L]} (a^{[L]} - y(x^{(k)}))$$

back prop.

for  $l = L-1$  down to 2

$$\delta^{[l]} = D^{[l]} (w^{[l+1]})^T \delta^{[l+1]}$$

end

for  $l = \{2 \text{ to } L\}$

gradient step

$$W^{[l]} \leftarrow W^{[l]} - \eta \delta^{[l]} (a^{[l-1]})^T$$

$$b^{[l]} \leftarrow b^{[l]} - \eta \delta^{[l]}$$

end

order doesn't matter

outer product  
with  $i, k$   
entry  
 $\delta_j^{[l]} a_k^{[l-1]}$

Efficiency is very dependent on step size/learning rate  $\eta$ .

Other Extensions Discussed in H&H

- convolutional neural networks
- avoiding overfitting using "validation data"
- other activation functions

$$\text{e.g. ReLU } \sigma(x) = \begin{cases} 0 & x \leq 0 \\ x & x \geq 0 \end{cases}$$

- other cost functions

e.g. softmax instead of  $\| \cdot \|_2^2$  ( $\approx \| \cdot \|_\infty$ )