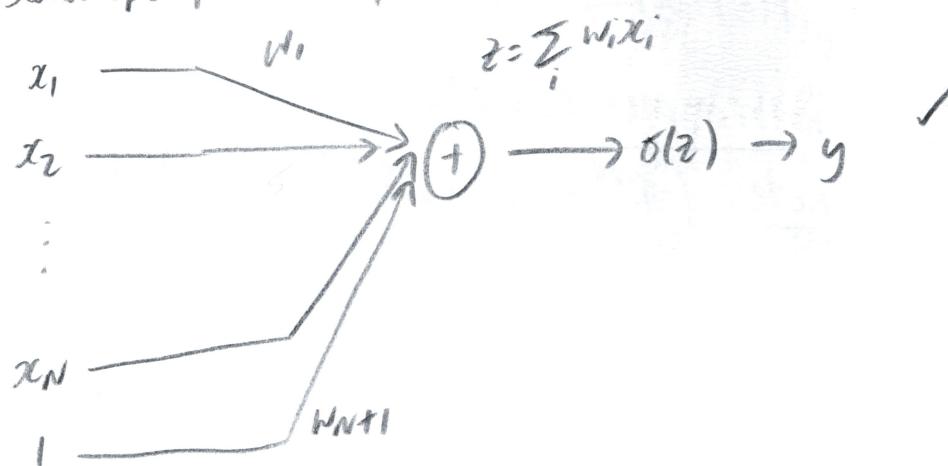


### 13-Training NN (review)

(\*) scalar specification of MLP with diff activations. (sigmoid)



$$\frac{\partial y}{\partial z} = \sigma'(z) \quad \frac{\partial y}{\partial w_i} = \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial w_i} = \sigma'(z)x_i \quad \frac{\partial y}{\partial x_i} = \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial x_i} = \sigma'(z)w_i$$

- okay  $\rightarrow$  happy with this

(\*) A little uncomfortable with the way the distinction between expected and empirical error; and the graphical presentation.

- whilst the presentation (BR) is consistent within presentation, trying to figure out discomfort

- in particular, the way in which the probabilistic interpretation is offered.

- the general presentation of generalisation error, training error in (expected) (empirical)

Bishop, Hastie, Salakhutdinov:-

- is that expectation of error tends to be taken over a joint distribution of inputs  $\underline{x}$  and target outputs  $t$

$$\text{i.e. } \mathbb{E}_{\underline{x}, t \sim p(\underline{x}, t)} [L] = \iint L(t, y(\underline{x})) p(\underline{x}, t) d\underline{x} dt$$

- AND a distinction made between classification and regression (via choice of loss function).

(\*) comfortable with derivatives; review gradient operator

## (\*) Gradient operator

for a scalar valued diff. function  $f: \mathbb{R}^n \mapsto \mathbb{R}$   
 of several variables

$$f(\underline{x}): \underline{x} \mapsto f(\underline{x})$$

$$\underline{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_N \end{pmatrix}$$

$$\nabla_{\underline{x}} f(\underline{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_N} \end{bmatrix}$$

- with respect to vector  $\underline{x}$
- note:  $f: \mathbb{R}^n \rightarrow \mathbb{R}$   $\underline{x} \in \mathbb{R}^{N \times 1}$   $\nabla_{\underline{x}} f(\underline{x}) \in \mathbb{R}^{N \times 1}$

## (\*) Note: understand that:-

- i) gradient of scalar fn of a vector points in rate of greatest increase of a function
  - ii) - if it is -ve; rate of fastest decrease of function
- } an incomplete specimen
- (\* see DL direct derivatives

(\*) ⓘ ⓘ → formal, mathematical justification for this; it's hand-wavy

- AS THIS IS NOT crucial; but would be helpful; leave it here;  
 find resources → add to Q/S list

- req will require additional machinery (vector calculus) - ⓘ

(\* directional derivative (\* see MIT OCW  
 or multivariable calculus OR "Deep learning")

(\* level sets (\* also clarifies the  
 approx. not. BR uses

(\*) later → gradient operator and level sets

i.e. for scalar valued functions of several variables

## (\*) Hessian matrix

- similar to gradient → matrix of 2nd

(setup)

$$\nabla^2 f(x_1, \dots, x_N) = \nabla^2 f(\underline{x}) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_N} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_N} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_N \partial x_1} & \frac{\partial^2 f}{\partial x_N \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_N^2} \end{bmatrix}$$

- (\*) can use properties of Hessian matrix i.e.  $S_n^{++}, S_n^+, S_n^-, S_n^-$   
 to identify local minima/maxima  
 and values of its eigenvalues (all of them)
- All eigenvalues of matrix  $\lambda_i > 0 \forall i \Rightarrow S_n^{++}$  (positive def)  $\rightarrow$  minimum
  - $\lambda_i < 0 \forall i \Rightarrow S_n^-$  (neg. def)  $\rightarrow$  maximum
- (\*) deep learning book really helps here. (goodfellow et al.)
- ch 9.3 - Gradient-based optimisation
  - 4.3.1 - Jacobian, Hessians
- (\*) the relation between expected error and empirical error and also empirical risk minimisation is given in DL (Goodfellow), ch 8 (unregularised) (empirical estimate of training exp. error) (8.1)
- $J(\theta) = \mathbb{E}_{(x,y)} \sim \hat{P}_{\text{data}} L(f(x;\theta), y)$
- (8.1) - Obj. function not training set  
 - we also have:-
- $$J^*(\theta) = \mathbb{E}_{(x,y)} \sim P_{\text{data}} L(f(x;\theta), y)$$
- this is the expected generalisation error or risk, where the expectation is taken over the data generating distribution  $P_{\text{data}}$  rather than just training set
- we do not know  $P_{\text{data}}(x,y)$ , so replace with  $\hat{P}_{\text{data}}(x,y)$  defined (empirical distn)  
 (true distn)  
 by training set.
- minimize  $\mathbb{E}_{x,y \sim \hat{P}_{\text{data}}(x,y)} L(f(x;\theta), y) = \frac{1}{m} \sum_{i=1}^m L(f(x^{(i)}; \theta), y^{(i)})$
- m - no. training examples

(米)

卷之三

10 of 10