

Defining NN bubbles →

1. Preprocess data →

A. Normalise

B. For images, only transcribe.

2. Choose NN Architecture →

3. Parameter initialization → depending on the NN.

4. 1st Step →

A. Disable regularisation.

B. Check if reasonable loss →

$$\mathcal{L} = \sum_i^{N=C} y_i \log p_i$$

2.3026...  $\times$   $\log \approx 2.3$ ; correct for 10 classes.

5. 2nd Step →

A. Crank up regularisation →  $1e^3 \Rightarrow$  loss ↑ 3.0685974

6. 3rd Step →

A. Small piece of data to overfit.

B. Regularisation = 0; num-epoch → loss  $\approx 0$ ; Accuracy = 100%

7. 4th Step → Find a learning rate that works

A.  $lr = 1e^{-6}$  → too small; loss changes.

B.  $lr = 1e^{-3}$  → too high; explodes.

Define system.

8. Hyperparameter search  $\rightarrow$

- \* Implement code to detect explosion and leak-act.
- \* By now rough idea of lr.

A. Coarse search  $\rightarrow$  5 epochs.

for count 100

$$reg = 10^{\wedge} \text{uniform}(-5, 5)$$

$$lr = 10^{\wedge} \text{uniform}(-3, -6)$$

$\rightarrow$  do search in log space  
because of low gradient  
noise.

learning decay = 0.9. } examples.

batch size = 100

look at Accuracy to define a finer search.

B. Fine search  $\rightarrow$

$$reg = 10^{\wedge} \text{uniform}(-5, 5)$$

$$lr = 10^{\wedge} \text{uniform}(-3, -6)$$

$$\text{uniform}(-4, 0)$$

$$\text{uniform}(-3, -4)$$

## NN Features →

### 1. Architecture →

- \* Deep ; Shallow
- \* Regularization
- \* Dropout.
- \* Weight initialization.
- \* Loss function.

### 2. Hyperparameters →

- \* Momentum →  $\beta_1 = 0.9$
  - \* RMSprop →  $\beta_2 = 0.999$
  - \* Learning decay.
- └─ Adam



$x \equiv$  Input; training set with  $q_x$  dimensions

$x \in (q_x, m) \in \mathbb{R} / q_x = \text{dimensions of set}; m = \text{training set length.}$

$y \equiv$  Result of binary classification

$y \in \mathbb{R} / y \in \{0, 1\} \Rightarrow Y = [y(1) \dots y(m)]$

$(x, y) / x \in \mathbb{R}^{(q_x, m)}; y \in \mathbb{R}^{(1, m)}$

\* Logistic Regression  $\rightarrow$

$\hat{y} = P(y=1|x)$ ; probability of  $y=1$  for  $x \rightarrow \hat{y}$  estimation.

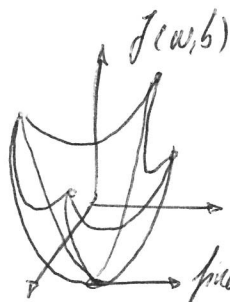
$$\hat{y} = \sigma(w^T \cdot x + b) = \sigma(z) = \frac{1}{1 + e^{-z}}$$

\* Loss function  $\rightarrow$  Error function.

$$\mathcal{L}(\hat{y}, y) = -(\log(p(y|x))) = -(y \log \hat{y} + (1-y) \log(1-\hat{y})) = \begin{cases} -\log \hat{y} & y=1 \\ (1-y) \log(1-\hat{y}) & y=0 \end{cases}$$

\* Cost function  $\rightarrow$

$$J(w, b) = \frac{1}{M} \sum_{i=1}^M \mathcal{L}(\hat{y}(i), y(i)) = -\frac{1}{M} \sum_{i=1}^M y(i) \log \hat{y}(i) + (1-y(i)) \log(1-\hat{y}(i))$$



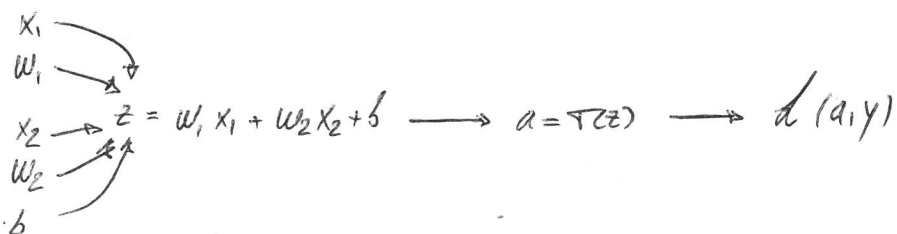
find  $w, b$  that minimizes the cost.

- Idea  $\rightarrow$  for training set, keep updating  $w, b$  value based on  $\partial J(w, b) / \partial w, \partial b$

$$w = w - \alpha \frac{\partial J(w, b)}{\partial w}$$

$$b = b - \alpha \frac{\partial J(w, b)}{\partial b}$$

\* Training set  $m=1 \rightarrow 2$  features.



$$d(a, y) = -[y \log(a) + (1-y) \log(1-a)] \rightarrow \frac{\partial d(a, y)}{\partial a} = -\frac{y}{a} + (y-1) \frac{1}{a-1}$$

$$a = S(z) = \frac{1}{1+e^{-z}} \rightarrow \frac{\partial a}{\partial z} = \frac{e^{-z}}{(1+e^{-z})^2} = a(1-a)$$

$$\frac{\partial d(a, y)}{\partial z} = a - y$$

- logistic regression for  $m$  examples  $\rightarrow M$  training set.

$J=0; \partial w_1=0; \partial w_2=0; \partial b=0 \rightarrow$  Initialization values.

for  $i=1$  to  $M$ :

$$z(i) = w^T x(i) + b \rightarrow (1 \times \eta_x) \times (\eta_x \times 1) + (1); \quad z(i) = \cancel{w_1^T x_1(i)} + \cancel{w_2^T x_2(i)} + b$$

$$a(i) = S(z(i)) \quad (1 \times 1)$$

$$J = -y [y(i) \log a(i) + (1-y(i)) \log(1-a(i))] \quad (1 \times 1)$$

$\uparrow$   
 $\eta_x = 2$ ; cross came from here.

$$dz(i) = a(i) - y(i) = a(i)(1-a(i)) \quad (1 \times 1)$$

$$dw_1 += x_1(i) dz(i)$$

$$dw_2 += x_2(i) dz(i)$$

$$db += dz(i) \quad (1 \times 1)$$

$\rightarrow$  number of features  $\rightarrow$  this will imply  $\eta_x = 2$ .

$$\hookrightarrow dw(i) += \frac{x(i)}{\eta_x \times 1} dz(i) \quad \frac{1 \times 1}{1 \times 1}$$

$$J = J/M; \quad dw_1 = dw_1/M; \quad dw_2 = dw_2/M; \quad db = db/M$$

- Vectorized logistic regression  $\rightarrow$

$$Z = W^T X + B \rightarrow (1 \times q_x)(q_x \times m) + (1 \times m) = (1 \times m)$$

$$A = \sigma(Z) \rightarrow (1 \times m)$$

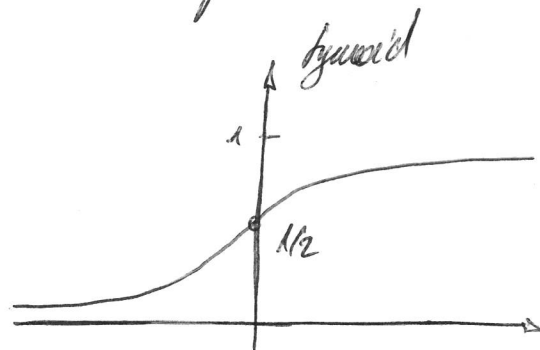
$$dz = A - Y \rightarrow (1 \times m)$$

$$dW = \frac{1}{M} X dz^T \rightarrow (q_x \times m) = (q_x \times m)(1 \times m)^T$$

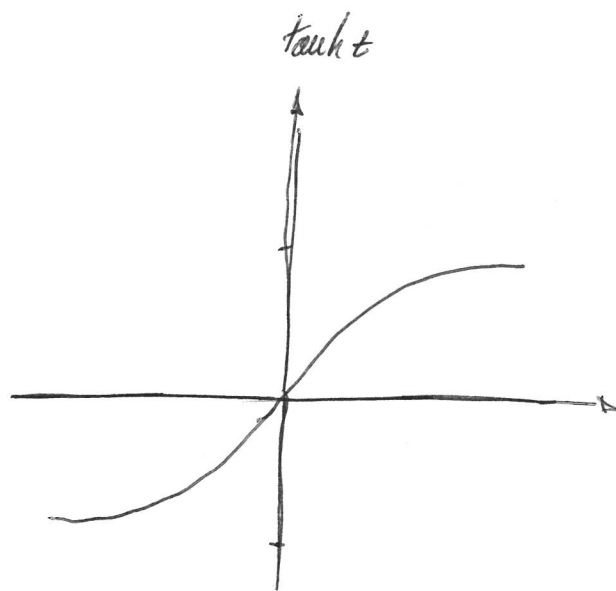
$$dB = \frac{1}{M} dz \rightarrow (1, 1)$$

\* Shallow Neural Networks  $\rightarrow$

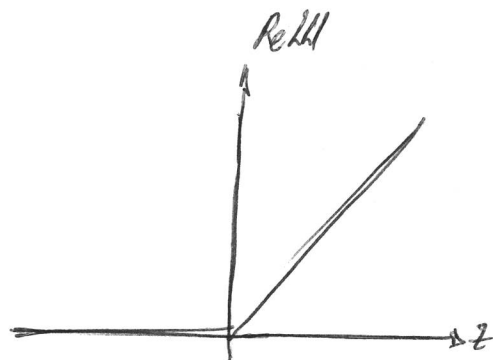
- Activation functions  $\rightarrow$



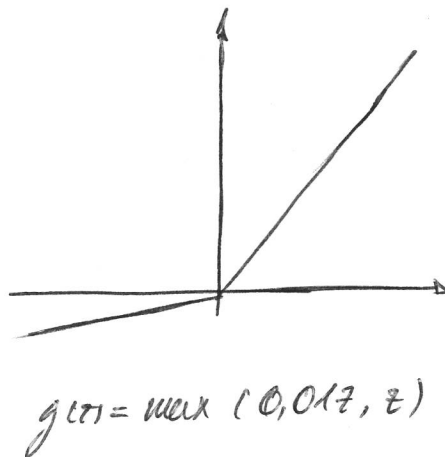
$$\sigma(z) = \frac{1}{1 + e^{-z}}; \sigma'(z) = \sigma(z)(1 - \sigma(z))$$



$$g(z) = \tanh(z) = \frac{e^z + e^{-z}}{e^z - e^{-z}}; g'(z) = 1 - \tanh^2(z)$$

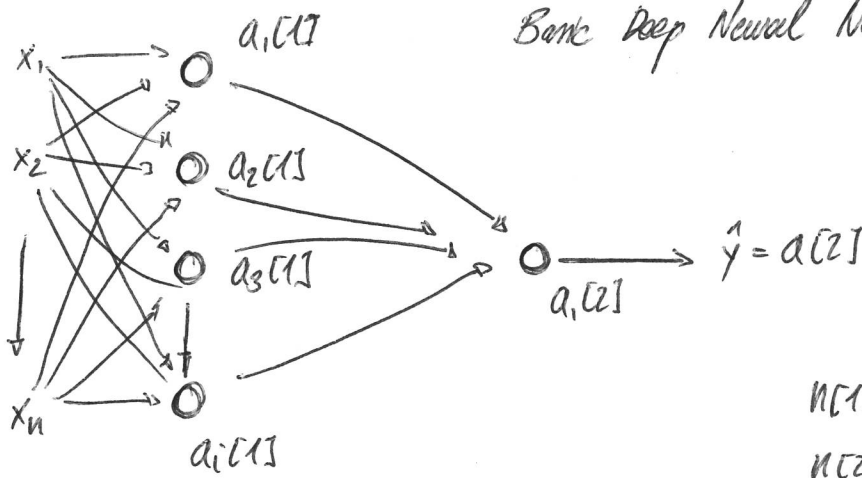


$$g(z) = \max(0, z); g'(z) = \begin{cases} 0 & z < 0 \\ 1 & z > 0 \\ \text{undef} & z = 0 \end{cases}$$



$$g'(z) = \begin{cases} 0.01 & z < 0 \\ 1 & z > 0 \\ \text{undef} & z = 0 \end{cases}$$

# Basic Deep Neural Networks



$$n[1] = I$$

$$n[2] = 1$$

layer [1] = hidden layer  
activation functions.

$$z_i[1] = w_i^T x + b \longrightarrow (1 \times m) = (\eta_x \times 1)^T (\eta_x \times m) + (1 \times m)$$

$$A_i[1] = g(z_i[1]) \longrightarrow (1 \times m) = (1 \times m)$$

$$A[1] = g(z[1]) \longrightarrow (I \times m) = (I \times m)$$

$$z[2] = w[2] A[1] + b[2] \longrightarrow (1 \times m) = (1 \times I)(I \times m) + (1 \times m)$$

$$A[2] = \sigma(z[2]) \longrightarrow (1 \times m)$$

$$dz[2] = A[2] - y \longrightarrow (1 \times m)$$

$$dw[2] = \frac{1}{M} dz[2] \cdot A[1]^T \longrightarrow (1 \times I) = (1 \times m)(I \times m)^T$$

$$db[2] = \frac{1}{M} \sum_{i=1}^M dz[2]$$

$$dz[1] = w[2]^T dz[2] * g'[1](z[1]) \longrightarrow (I \times m) = (1 \times I)^T (1 \times m) * (I \times m)$$

$$dw[1] = dz[1] \cdot x^T \longrightarrow (\eta_x \times I) = (I \times m)(\eta_x \times m)^T$$

$$db[1] = dz[1]$$



\* Initialization gradient descent  $\rightarrow$

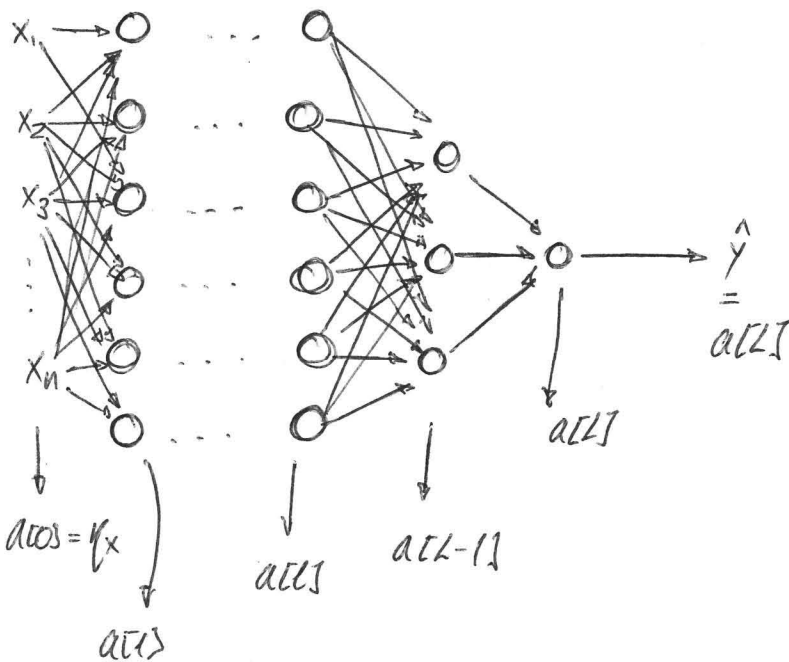
$$\frac{\partial J(a^{(2)}, y)}{\partial w_j^{(1)}} = \frac{1}{M} \sum_{i=1}^M \frac{\partial L(a^{(2)}, y)}{\partial w_j^{(1)}}$$

$$\frac{1}{M} \sum_{i=1}^M \frac{\partial L(a^{(2)}, y)}{\partial a^{(2)}} \cdot \frac{\partial a^{(2)}}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial a_i^{(1)}} \cdot \frac{\partial a_i^{(1)}}{\partial z_i^{(1)}} \cdot \frac{\partial z_i^{(1)}}{\partial w_j^{(1)}} =$$

$$\frac{1}{M} \sum_{i=1}^M (a^{(2)} - y) [w_j^{(2)} \cdot g'_i(z_i^{(1)})] \cdot x$$

$$dW^{(1)} = \begin{bmatrix} x_1 w_1^{(2)} g'_1(z_1^{(1)}) & \dots & x_1 w_n^{(2)} g'_n(z_n^{(1)}) \\ \vdots & & \vdots \end{bmatrix}$$

- Deep 6-layer Neural Network  $\rightarrow$



\* Forward Propagation  $\rightarrow$

$$z^{(L)} = W^{(L)} \cdot A^{(L-1)} + b^{(L)}$$

$$a^{(L)} = g^{(L)}(z^{(L)})$$

$$W \in \mathbb{R}(n^{(L)}, n^{(L-1)}); b \in \mathbb{R}(n^{(L)}, 1)$$

$$A \in \mathbb{R}(n^{(L)}, n^{(L-1)})$$

\* Backward propagation  $\rightarrow$

$$d z^{(L)} = d A^{(L)} * g^{(L)}(z^{(L)})$$

$$d A^{(L-1)} = W^{(L)T} d z^{(L)}$$

$$d W^{(L)} = \frac{1}{M} d z^{(L)} A^{(L-1)T}$$

$$d b^{(L)} = \frac{1}{M} \sum_{i=1}^M d z^{(L)}_i$$

$$\begin{aligned}
 X &\longrightarrow z_{CLS} = W_{CLS} A_{CLS} + B_{CLS} && \xrightarrow{A_{CLS}} \dots \xrightarrow{A_{CLS-1}} z_{CLS} = W_{CLS} A_{CLS-1} + B_{CLS} && \xrightarrow{A_{CLS}} \\
 &A_{CLS} = g_{CLS}(z_{CLS}) && && A_{CLS} = g_{CLS}(z_{CLS})
 \end{aligned}$$

$$\begin{aligned}
 &\downarrow A_{CLS}, W_{CLS} \\
 dz_{CLS} &= dA_{CLS} * g'_{CLS}(z_{CLS})
 \end{aligned}$$

$$dW_{CLS} = \frac{1}{M} dz_{CLS} A_{CLS}^T \longleftarrow \dots \longleftarrow \begin{matrix} dA_{CLS-1} \\ dA_{CLS} \end{matrix}$$

$$dB_{CLS} = \frac{1}{M} \sum_{i=1}^M dz_{CLS}$$

$$dA_{CLS} = W_{CLS}^T dz_{CLS}$$

$$\downarrow \\
 dW_{CLS}, dB_{CLS}$$

$$\begin{aligned}
 dz_{CLS} &= dA_{CLS} * g'_{CLS}(z_{CLS}) \\
 dW_{CLS} &= \frac{1}{M} dz_{CLS} A_{CLS-1}^T \longleftarrow \dots \longleftarrow \begin{matrix} dA_{CLS-1} \\ dA_{CLS} \end{matrix} \\
 dB_{CLS} &= \frac{1}{M} \sum_{i=1}^M dz_{CLS}
 \end{aligned}$$

$$dA_{CLS-1} = W_{CLS}^T dz_{CLS}$$

$$\downarrow \\
 dW_{CLS}, dB_{CLS}$$

$$\begin{aligned}
 A_{CLS-1} &\longrightarrow z_{CLS} = W_{CLS} A_{CLS-1} + B_{CLS} && \longrightarrow A_{CLS} = \hat{y} \\
 &A_{CLS} = g_{CLS}(z_{CLS})
 \end{aligned}$$

$$\downarrow A_{CLS-1}, W_{CLS}$$

$$dz_{CLS} = dA_{CLS} * g'_{CLS}(z_{CLS})$$

$$dW_{CLS} = \frac{1}{M} dz_{CLS} A_{CLS-1}^T \longleftarrow \dots \longleftarrow dA_{CLS} = -\frac{y}{A_{CLS}} + \frac{(1-y)}{(1-A_{CLS})}$$

$$\begin{aligned}
 &\longleftarrow \dots \longleftarrow dA_{CLS-1} \\
 dB_{CLS} &= \frac{1}{M} \sum_{i=1}^M dz_{CLS}
 \end{aligned}$$

$$dA_{CLS-1} = W_{CLS}^T dz_{CLS}$$

## Backpropagation notes

$$A = g(z) ; A, z \in \mathbb{R}^{(n_A, m)}$$

$$z = W A_{prev} ; A_{prev} \in \mathbb{R}^{(n_x, m)} ; W \in \mathbb{R}^{(n_A, n_x)}$$

$$\frac{\partial z_{ke}}{\partial w_{ij}} = \begin{cases} 0 & \text{if } k \neq i \\ A_{prev e j} & \text{if } k = i \end{cases} ; \quad \frac{\partial A_{ke}}{\partial z_{ij}} = g'(z_{ke}) \delta_{k,i} \delta_{e,j}$$

$$\begin{aligned} \frac{\partial A_{ke}}{\partial w_{ij}} &= \sum_{k=1}^{n_A} \sum_{f=1}^m \frac{\partial A_{ke}}{\partial z_{kf}} \cdot \frac{\partial z_{kf}}{\partial w_{ij}} = \sum_{k=1}^{n_A} \sum_{f=1}^m g'(z_{ke}) \cdot \delta_{k,e} \cdot \delta_{e,f} \cdot A_{prev f j} \delta_{k,i} \\ &= g'(z_{ke}) A_{prev e j} \delta_{k,i} \end{aligned}$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w_{ij}} &= \sum_{k=1}^{n_A} \sum_{e=1}^M \frac{\partial \mathcal{L}}{\partial A_{ke}} \frac{\partial A_{ke}}{\partial w_{ij}} = \sum_{k=1}^{n_A} \sum_{e=1}^M \frac{\partial \mathcal{L}}{\partial A_{ke}} g'(z_{ke}) \cdot A_{prev e j} \delta_{k,i} = \\ &= \sum_{e=1}^m \frac{\partial \mathcal{L}}{\partial A_{ie}} \cdot g'(z_{ie}) \cdot A_{e j} = \left( \frac{\partial \mathcal{L}}{\partial A_{i1}} g'(z_{i1}) A_{1j} + \dots + \frac{\partial \mathcal{L}}{\partial A_{i n_A}} g'(z_{i n_A}) A_{n_A j} \right) \end{aligned}$$

\* Vectorised implementation  $\rightarrow$

$$dW = (dA \cdot g'(z))_i A_{prev}^T$$

element-wise multiplication.  $\nearrow$  dot product

$$\frac{\partial \mathcal{L}}{\partial A_{prev,ij}} = W_{ki} \delta_{l,j} \quad ; \quad \frac{\partial A_{xc}}{\partial A_{prev,ij}} = \sum_{c=1}^{n_A} \sum_{f=1}^m \frac{\partial A_{xc}}{\partial A_{cf}} \cdot \frac{\partial A_{cf}}{\partial A_{prev,ij}} = g'(z_{xc}) \cdot W_{ki} \cdot \delta_{l,j}$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial A_{prev,ij}} &= \sum_{c=1}^{n_A} \sum_{f=1}^m \frac{\partial \mathcal{L}}{\partial A_{cf}} \cdot \frac{\partial A_{cf}}{\partial A_{prev,ij}} = \sum_{c=1}^{n_A} \sum_{f=1}^m \frac{\partial \mathcal{L}}{\partial A_{cf}} \cdot g(z_{cf}) W_{ki} \delta_{f,j} \\ &= \sum_{c=1}^{n_A} \frac{\partial \mathcal{L}}{\partial A_{cj}} g'(z_{cj}) W_{ki} \end{aligned}$$

\* Vectorised implementation.  $\longrightarrow$

$$dA_{prev} = W^T \cdot (dAL * g'(z))$$