

# Αναγωγή Προτύπων

## 2<sup>η</sup> Σειρά Αναλυτικών Ασκήσεων

Δημήτριος Κοκκίνης

03118896

dimkok00@gmail.com

### Άσκηση 2.1

1. Έχουμε το δίκτυο  $f(x) = f^L(f^{L-1}(f^{L-2}(\dots f^1(x))))$  όπου  $1, \dots, L-2, L-1, L$  οι ενότητες αυτού και  $f^l = \sigma(w_l x + b_l)$  η συνάρτηση που παρακτυπεί κάθε layer  $l$ , με παραμέτρους  $w_l$  και  $b_l$ . Θεωρώντας γραμμική συνάρτηση ενεργοποίησης  $\sigma(x) = x$ , τότε το δίκτυο γραφίζεται σαν μορφή:

$$\begin{aligned} f(x) &= f^L(f^{L-1}(f^{L-2}(\dots f^4(f^3(f^2(f^1(x))))))) = \\ &= f^L(f^{L-1}(f^{L-2}(\dots f^4(f^3(f^2(\sigma(w_1 x + b_1))))))) = \\ &= f^L(f^{L-1}(f^{L-2}(\dots f^4(f^3(\sigma(w_2(w_1 x + b_1) + b_2)))))) = \\ &= \dots = f^L(f^{L-1}(f^{L-2}(\dots f^4(w_3 w_2 w_1 x + w_3 w_2 b_1 + w_3 b_2 + b_3)))) \\ &= \dots = w_{L-1} w_{L-2} \dots w_1 (x) + w_{L-1} w_{L-2} \dots w_2 b_1 x + w_{L-1} w_{L-2} \dots w_3 b_2 x + \dots + w_{L-1} b_{L-2} x + b_{L-1} \\ &= (w_{L-1} w_{L-2} \dots w_1 + w_{L-1} w_{L-2} \dots w_2 b_1 + w_{L-1} w_{L-2} \dots w_3 b_2 + \dots + w_{L-1} b_{L-2}) x + b_{L-1} \end{aligned}$$

Οπότε τώρα, η συνάρτηση  $f(x)$  είναι σαν μορφή  $w x + b$  και επειδή  $\sigma(x) = x$  άρα  $\sigma(w x + b) = w x + b$ . Όμως η  $\sigma(w x + b)$  εκφράζει ένα single layer network  $g(x) = w x + b$ .

Είναι  $\|f(x) - g(x)\| < \epsilon$ ,  $\epsilon > 0$  άρα από θεωρήμα καθολικής προσέγγισης η  $f(x)$  είναι αντιστοιχη με την  $g(x) = w x + b$  δηλαδή  $f(x) \approx g(x) = w x + b$ .

2. Από τα δεδομένα διαπιστώνουμε ότι έχουμε 2 κόμβους,  $x_1, x_2$ , στο input layer, έχουμε ένα deep layer,  $a$ , με 4 κόμβους  $a_1, a_2, a_3, a_4$ .

Για το ερώτημα αυτό ισχύει ότι:

$$\frac{\partial a}{\partial x} = \begin{bmatrix} \frac{\partial a_1}{\partial x_1} & \frac{\partial a_1}{\partial x_2} \\ \frac{\partial a_2}{\partial x_1} & \frac{\partial a_2}{\partial x_2} \\ \frac{\partial a_3}{\partial x_1} & \frac{\partial a_3}{\partial x_2} \\ \frac{\partial a_4}{\partial x_1} & \frac{\partial a_4}{\partial x_2} \end{bmatrix} = W_1 = \begin{bmatrix} 2 & 2 \\ -2 & -2 \\ 2 & -2 \\ -2 & 2 \end{bmatrix} = J^{(1)} \quad \text{ταυτοποίηση!}$$

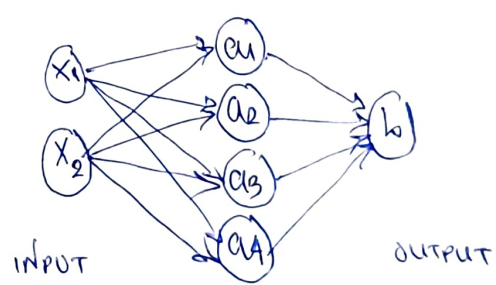
αλλά είναι  $a = W_1 x + b_1 = \begin{bmatrix} 2 & 2 \\ -2 & -2 \\ 2 & -2 \\ -2 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 & 0 \end{bmatrix}^T$

$$= \begin{bmatrix} 2(x_1 + x_2) \\ -2(x_1 + x_2) \\ 2(x_1 - x_2) \\ -2(x_1 - x_2) \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \end{bmatrix}$$

Έχουμε ένα output layer,  $b$ , με ένα κόμβο  $b$  για το οποίο ισχύει:

$$\frac{\partial b}{\partial x} = \begin{bmatrix} \frac{\partial b}{\partial a_1} & \frac{\partial b}{\partial a_2} & \frac{\partial b}{\partial a_3} & \frac{\partial b}{\partial a_4} \end{bmatrix} = W_2 = \begin{bmatrix} \mu & \mu & -\mu & -\mu \end{bmatrix} = J^{(2)} \quad \text{ταυτοποίηση!}$$

Επομένως το δίκτυο είναι το εξής:



Για το output layer ισχύει:

$$\begin{aligned} \tilde{f}(x_1, x_2) &= b = W_2 \sigma(W_1 x + b_1) = W_2 \sigma(a) = \begin{bmatrix} \mu & \mu & -\mu & -\mu \end{bmatrix} \begin{bmatrix} \sigma(a_1) \\ \sigma(a_2) \\ \sigma(a_3) \\ \sigma(a_4) \end{bmatrix} \\ &= \mu \sigma(a_1) + \mu \sigma(a_2) - \mu \sigma(a_3) - \mu \sigma(a_4) \\ &= \mu \sigma(2(x_1 + x_2)) + \mu \sigma(-2(x_1 + x_2)) - \mu \sigma(2(x_1 - x_2)) - \mu \sigma(-2(x_1 - x_2)) \end{aligned}$$

3. Εκπεράτουμε τη  $\tilde{f}(x_1, x_2)$  ως quadratic γύρω από το  $(0,0)$ :

$$\begin{aligned} \tilde{f}(x_1, x_2) &= \sigma(x_1, x_2) = \tilde{f}(0,0) + \tilde{f}_{x_1}(0,0)x_1 + \tilde{f}_{x_2}(0,0)x_2 + \frac{\tilde{f}_{x_1 x_1}(0,0)}{2} x_1^2 + \frac{\tilde{f}_{x_2 x_2}(0,0)}{2} x_2^2 \\ &\quad + \frac{\tilde{f}_{x_1 x_2}(0,0)}{2} x_1 x_2 \end{aligned}$$

$$\bullet \tilde{f}(0,0) = \mu \delta(0) + \mu \delta(0) - \mu \delta(0) - \mu \delta(0) = 0$$

$$\bullet \tilde{f}_{x_1}(0,0) = \mu \lambda \delta'(0) - \mu \lambda \delta'(0) - \mu \lambda \delta'(0) + \mu \lambda \delta'(0) = 0$$

$$\bullet \tilde{f}_{x_2}(0,0) = \mu \lambda \delta'(0) - \mu \lambda \delta'(0) + \mu \lambda \delta'(0) - \mu \lambda \delta'(0) = 0$$

$$\bullet \tilde{f}_{x_1 x_2}(0,0) = \frac{\partial \tilde{f}_{x_1}(0,0)}{\partial x_2} = \dots = \mu \lambda^2 \ddot{\delta}(0) + \mu \lambda^2 \ddot{\delta}(0) + \mu \lambda^2 \ddot{\delta}(0) + \mu \lambda^2 \ddot{\delta}(0) \\ = 4 \mu \lambda^2 \ddot{\delta}(0)$$

$$\bullet \tilde{f}_{x_1 x_1}(0,0) = \frac{\partial \tilde{f}_{x_1}(0,0)}{\partial x_1} = \dots = \mu \lambda^2 \ddot{\delta}(0) + \mu \lambda^2 \ddot{\delta}(0) - \mu \lambda^2 \ddot{\delta}(0) - \mu \lambda^2 \ddot{\delta}(0) = 0$$

$$\bullet \tilde{f}_{x_2 x_2}(0,0) = \frac{\partial \tilde{f}_{x_2}(0,0)}{\partial x_2} = \dots = \mu \lambda^2 \ddot{\delta}(0) + \mu \lambda^2 \ddot{\delta}(0) - \mu \lambda^2 \ddot{\delta}(0) - \mu \lambda^2 \ddot{\delta}(0) = 0$$

$$\text{οότε } \tilde{f}(x_1, x_2) = 4 \mu \lambda^2 \ddot{\delta}(0) x_1 x_2 \quad (1) \text{ και } \lim_{\lambda \rightarrow 0} \tilde{f}(x) = \lim_{\lambda \rightarrow 0} \tilde{f}(x_1, x_2)$$

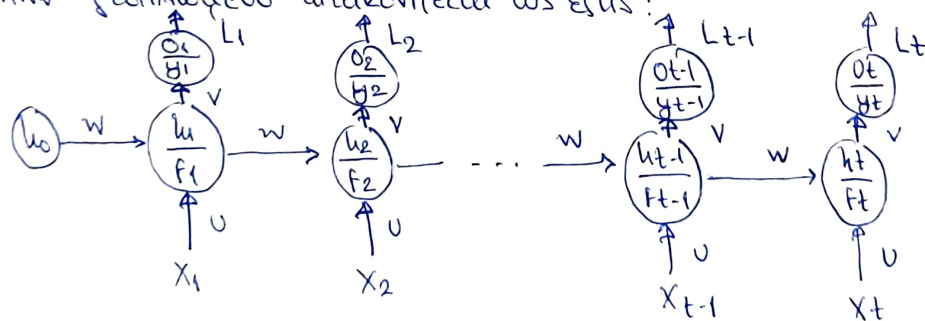
$$\stackrel{(1)}{=} \lim_{\lambda \rightarrow 0} 4 \mu \lambda^2 \ddot{\delta}(0) x_1 x_2$$

$$\underline{\underline{\ddot{\delta}(0) = 1/4\lambda^2 \mu}} \quad \lim_{\lambda \rightarrow 0} x_1 x_2 = f(x)$$

Επομένως για  $\lambda \rightarrow 0$  ισχύει ότι  $\tilde{f}(x) \rightarrow f(x)$

## Άσκηση 2.2

Το RNN περιγράφετο απεικονίζεται ως εξής:



1. Σύμφωνα με το κανόνα της αλυσίδας η παραγώγος  $\frac{\partial L_t}{\partial v}$  μπορεί να γραφεί ως

$$\text{εξής: } \frac{\partial L_t}{\partial v} = \frac{\partial L_t}{\partial o_t} \cdot \frac{\partial o_t}{\partial y_t} \cdot \frac{\partial y_t}{\partial v} \quad (1)$$

Είναι  $\frac{\partial y_t}{\partial v} = h_t$

Επομένως η (1) γράφεται:  $\boxed{\frac{\partial L_t}{\partial v} = \frac{\partial L_t}{\partial o_t} \cdot \frac{\partial o_t}{\partial y_t} \cdot h_t}$

2. Η  $\frac{\partial L_t}{\partial w}$  γράφεται:  $\frac{\partial L_t}{\partial w} = \frac{\partial L_t}{\partial o_t} \cdot \frac{\partial o_t}{\partial y_t} \cdot \frac{\partial y_t}{\partial h_t} \cdot \frac{\partial h_t}{\partial h_e} \cdot \frac{\partial h_e}{\partial w} = \frac{\partial L_t}{\partial o_t} \cdot \frac{\partial o_t}{\partial h_t} \cdot \frac{\partial h_t}{\partial h_e} \cdot \frac{\partial h_e}{\partial w} \quad (2)$

ο όρος  $\frac{\partial h_t}{\partial h_e}$  είναι η παράγωγος του hidden state τη χρονική στιγμή  $t$  αναφορικά με το hidden state τη χρονική στιγμή  $e$ .

Τον όρο αυτό υποθέτουμε να τον γράψουμε ως:

$$\frac{\partial h_t}{\partial h_{t-1}} \cdot \frac{\partial h_t}{\partial h_e} = \frac{\partial h_t}{\partial h_{t-1}} \cdot \frac{\partial h_{t-1}}{\partial h_{t-2}} \cdots \frac{\partial h_{t+1}}{\partial h_e} = \prod_{i=t+1}^t \frac{\partial h_i}{\partial h_{i-1}} \quad (3)$$

και λόγω της σχέσης  $h_t = f(w h_{t-1} + v x_t)$  η (3) γράφεται ως εξής:

$$\frac{\partial h_t}{\partial h_e} = \prod_{i=t+1}^t \frac{\partial h_i}{\partial h_{i-1}} = \prod_{i=t+1}^t w^T \text{diag}[f'(h_{i-1})] \quad (4)$$

Επομένως η (2) από (4), (5) γράφεται:  $\boxed{\frac{\partial L_t}{\partial w} = \frac{\partial L_t}{\partial o_t} \cdot \frac{\partial o_t}{\partial h_t} \cdot \left( \prod_{i=t+1}^t w^T \text{diag}[f'(h_{i-1})] \right) \cdot \frac{\partial h_e}{\partial w}}$

3. Θεωρούμε ότι η activation function  $f$  είναι η ταυτοτική συνάρτηση.

Τότε είναι:  $\frac{\partial h_t}{\partial h_e} = \prod_{i=t+1}^t \frac{\partial h_i}{\partial h_{i-1}} = (w^T)^{t-1} \stackrel{t-e=k}{=} (w^T)^k$

Έστω  $\lambda_1, \lambda_2, \dots, \lambda_n$  οι ιδιοτιμές του  $w$  με  $|\lambda_1| > |\lambda_2| > \dots > |\lambda_n|$  και τα αντίστοιχα ιδιοδιανύσματα  $v_1, v_2, \dots, v_n$  να σχηματίζουν μια διανυσματική βάση.

ο  $w$  είναι τετραγωνικός, διαγωνοποιήσιμο και μπορούμε να χρησιμοποιήσουμε τη σχέση  $v_i^T (w^T)^k = \lambda_i^k v_i^T$ . Μπορούμε επίσης να γράψουμε το διάνυσμα

γραμμής  $\frac{\partial L_t}{\partial h_t}$  χρησιμοποιώντας τη βάση  $v_1, v_2, \dots, v_n$ :  $\frac{\partial L_t}{\partial h_t} = \sum_{i=1}^N c_i v_i^T$

Επιλέγουμε  $j$  τέτοιο ώστε  $c_j \neq 0$  και για κάθε  $j' < j$  είναι  $c_{j'} = 0$ .

Έχουμε:  $\frac{\partial L_t}{\partial h_t} \frac{\partial h_t}{\partial h_e} = c_j \lambda_j^k v_j^T + \lambda_j^k \sum_{i=j+1}^n c_i \frac{\partial \lambda_i^k}{\partial \lambda_j^k} v_i^T \quad (6)$

Για  $i > j$  είναι  $|z_i| < |z_j| \Rightarrow \left| \frac{z_i}{z_j} \right| < 1$  επομένως  $\lim_{l \rightarrow \infty} \left| \frac{z_i}{z_j} \right|^l = 0$  οπότε

η (6) απλοποιείται ως:

$$\frac{\partial L_t}{\partial h_t} \cdot \frac{\partial h_t}{\partial h_k} \approx c_j z_j^k V_j^T$$

Αν  $|z_j| > 1$  τότε  $\frac{\partial h_t}{\partial h_k} = z_j^k$  αυξάνεται εκθετικά καθώς η διεύθυνση των  $V_j$  και επομένως το gradient  $\frac{\partial L_t}{\partial w}$  να είναι ίσο με  $\frac{\partial L_t}{\partial t} \cdot \frac{\partial t}{\partial h_t} \cdot \frac{\partial h_t}{\partial h_k} \cdot \frac{\partial h_k}{\partial w}$  μεγαίνει στο άπειρο. (exploding gradient),

Αν  $|z_j| < 1$  τότε  $\frac{\partial h_t}{\partial h_k} = z_j^k$  μειώνεται εκθετικά καθώς η διεύθυνση των  $V_j$  και επομένως το gradient τείνει να εξαφανιστεί (vanishing gradient),



Άσκηση 2.4

$$\left. \begin{aligned} \mu_1 &= \vec{w}^T \vec{\mu}_1, \quad \bar{z}_1 = \sum_i (\vec{x}_i - \vec{\mu}_1)(\vec{x}_i - \vec{\mu}_1)^T \\ \mu_2 &= \vec{w}^T \vec{\mu}_2, \quad \bar{z}_2 = \sum_i (\vec{x}_i - \vec{\mu}_2)(\vec{x}_i - \vec{\mu}_2)^T \end{aligned} \right\} \textcircled{A}$$

$$\Rightarrow \left. \begin{aligned} \sigma_1^2 &= \sum_i (y_i - \mu_1)^2 = \sum_i (\vec{w}^T \vec{x}_i - \vec{w}^T \vec{\mu}_1)^2 \\ \sigma_2^2 &= \sum_i (y_i - \mu_2)^2 = \sum_i (\vec{w}^T \vec{x}_i - \vec{w}^T \vec{\mu}_2)^2 \end{aligned} \right\} \Leftrightarrow \begin{aligned} \sigma_1^2 &= \sum_i \vec{w}^T (\vec{x}_i - \vec{\mu}_1)(\vec{x}_i - \vec{\mu}_1)^T \vec{w} \\ \sigma_2^2 &= \sum_i \vec{w}^T (\vec{x}_i - \vec{\mu}_2)(\vec{x}_i - \vec{\mu}_2)^T \vec{w} \end{aligned}$$

$$\Leftrightarrow \left. \begin{aligned} \sigma_1^2 &= \vec{w}^T \sum_i (\vec{x}_i - \vec{\mu}_1)(\vec{x}_i - \vec{\mu}_1)^T \vec{w} \\ \sigma_2^2 &= \vec{w}^T \sum_i (\vec{x}_i - \vec{\mu}_2)(\vec{x}_i - \vec{\mu}_2)^T \vec{w} \end{aligned} \right\} \Leftrightarrow \begin{aligned} \sigma_1^2 &= \vec{w}^T \bar{z}_1 \vec{w} \\ \sigma_2^2 &= \vec{w}^T \bar{z}_2 \vec{w} \end{aligned}$$

Έτσι η Fisher LDA των παραπάνω καθορίζεται από

$$J_F(\vec{w}) = \frac{(\mu_1 - \mu_2)^2}{\sigma_1^2 + \sigma_2^2} = \frac{S_B^y}{S_W^y} = \frac{\vec{w}^T S_B \vec{w}}{\vec{w}^T S_W \vec{w}}$$

καθώς  $S_B^y = (\mu_1 - \mu_2)^2 = (\vec{w}^T \vec{\mu}_1 - \vec{w}^T \vec{\mu}_2)^2 = \vec{w}^T (\vec{\mu}_1 - \vec{\mu}_2)(\vec{\mu}_1 - \vec{\mu}_2)^T \vec{w} = \vec{w}^T S_B \vec{w}$

$$S_W^y = \sigma_1^2 + \sigma_2^2 = (\vec{w}^T \bar{z}_1 \vec{w}) + (\vec{w}^T \bar{z}_2 \vec{w}) = \vec{w}^T (\bar{z}_1 + \bar{z}_2) \vec{w} = \vec{w}^T S_W \vec{w}$$

Αν  $S_W$  αντιστρέψιμο:

$$\frac{dJ_F(\vec{w})}{d\vec{w}} = \frac{d}{d\vec{w}} \left( \frac{\vec{w}^T S_B \vec{w}}{\vec{w}^T S_W \vec{w}} \right) = \frac{(\vec{w}^T S_W \vec{w}) 2 S_B \vec{w} - (\vec{w}^T S_B \vec{w}) 2 S_W \vec{w}}{(\vec{w}^T S_W \vec{w})^2} = 0$$

$$\Leftrightarrow 2(\vec{w}^T S_W \vec{w}) S_B \vec{w} - 2(\vec{w}^T S_B \vec{w}) S_W \vec{w} = 0 \Leftrightarrow \vec{w}^T S_W \vec{w} S_B \vec{w} = \vec{w}^T S_B \vec{w} S_W \vec{w} \Leftrightarrow$$

$$\Leftrightarrow S_B \vec{w} = J_1(\vec{w}) S_W \vec{w} \Leftrightarrow (S_W^{-1} S_B) \vec{w} = J_1(\vec{w}) \vec{w} \begin{matrix} \text{ενδιαφέρουσες} \\ \text{ιδιοτιμές} \end{matrix} \Leftrightarrow \vec{w}^* = S_W^{-1} (\vec{\mu}_1 - \vec{\mu}_2)$$

$$\Leftrightarrow \boxed{\vec{w}^* = (\bar{z}_1 + \bar{z}_2)^{-1} (\vec{\mu}_1 - \vec{\mu}_2)} \quad \text{δίνει αν } S_W \text{ αντιστρέψιμο}$$

$$S_W^{-1} S_B \vec{w} = \lambda \vec{w} \Leftrightarrow S_W^{-1} (\vec{\mu}_1 - \vec{\mu}_2) (\vec{\mu}_1 - \vec{\mu}_2)^T \vec{w} = \lambda \vec{w} \Leftrightarrow S_W^{-1} (\vec{\mu}_1 - \vec{\mu}_2) = \frac{1}{\lambda} \vec{w} \Leftrightarrow$$

$$\Leftrightarrow \vec{w} = \frac{a}{\lambda} (\bar{z}_1 + \bar{z}_2)^{-1} (\vec{\mu}_1 - \vec{\mu}_2) \Rightarrow \vec{w} = (\bar{z}_1 + \bar{z}_2)^{-1} (\vec{\mu}_1 - \vec{\mu}_2) \quad \text{γιατί δεν μας ενδιαφέρει το } \|\vec{w}\|^2$$

# Άσκηση 2.5

$$1. \text{ Υπολογίζουμε } P(G=1) = P(G=1|B=1)P(B=1) + P(G=1|B=0)P(B=0) \text{ όπου } P(G=1|B=1) \\ = P(G=1|B=1, F=1)P(F=1) + P(G=1|B=1, F=0)P(F=0) = 0.95 \cdot 0.8 + 0.3 \cdot 0.2 = 0.82$$

$$P(G=1|B=0) = \dots = 0.24 \text{ οπότε } P(G=1) = 0.791$$

$$\square P(D=0) = P(D=0|G=1)P(G=1) + P(D=0|G=0)P(G=0) = 0.2 \cdot 0.791 + 0.8 \cdot 0.209 = 0.3254$$

$$\square P(G=0|F=0) = P(G=0|B=0, F=0)P(B=0|F=0) + P(G=0|B=1, F=0)P(B=1|F=0) = \\ = 0.8 \cdot 0.05 + 0.7 \cdot 0.95 = 0.705$$

$$\square P(D=0|F=0) = 0.8 \cdot 0.705 + 0.2 \cdot 0.795 = 0.623 \text{ άρα } P(F=0|D=0) = \dots = 0.383$$

2.

$$P(F=0|D=0, B=0) = \frac{P(F=0, D=0, B=0)}{P(D=0, B=0)} = \frac{P(D=0|F=0, B=0)P(F=0, B=0)}{P(D=0, B=0)} \\ = P(D=0|F=0, B=0) \frac{P(F=0|B=0)}{P(D=0|B=0)}$$

$$\square P(D=0|F=0, B=0) = P(D=0|G=0)P(G=0|F=0, B=0) + P(D=0|G=1)P(G=1|F=0) \\ = 0.8 \cdot 0.8 + 0.2 \cdot 0.2 = 0.68$$

$$\square P(D=0|B=0) = P(D=0|G=0)P(G=0|B=0) + P(D=0|G=1)P(G=1|B=0)$$

$$\square P(G=1|B=0) = P(G=1|B=0, F=0)P(F=0) + P(G=1|B=0, F=1)P(F=1) \\ = 0.2 \cdot 0.2 + 0.25 \cdot 0.8 = 0.24$$

$$\text{οπότε } P(D=0|B=0) = 0.8 \cdot 0.76 + 0.2 \cdot 0.24 = 0.656$$

$$P(F=0|D=0, B=0) = 0.207$$

Ανταδύ!  $P(F=0|D=0, B=0) < P(F=0|D=0) = 0.383$  Διότι στην  $P(F=0|D=0)$  έχουμε λαβει υπόψη τη πιθανότητα  $B=0$  ή  $B=1$  καί' με τη παρατήρηση των  $G$