# Review

$d_1$

$d_2$

sigmoid

$\sigma$

$W_1, b_1$

$W_2, b_2$

$W_3, b_3$

Dense_1

$P \in [0,1]$

$X \in \mathbb{R}^d$

$z^1$

$z^2$

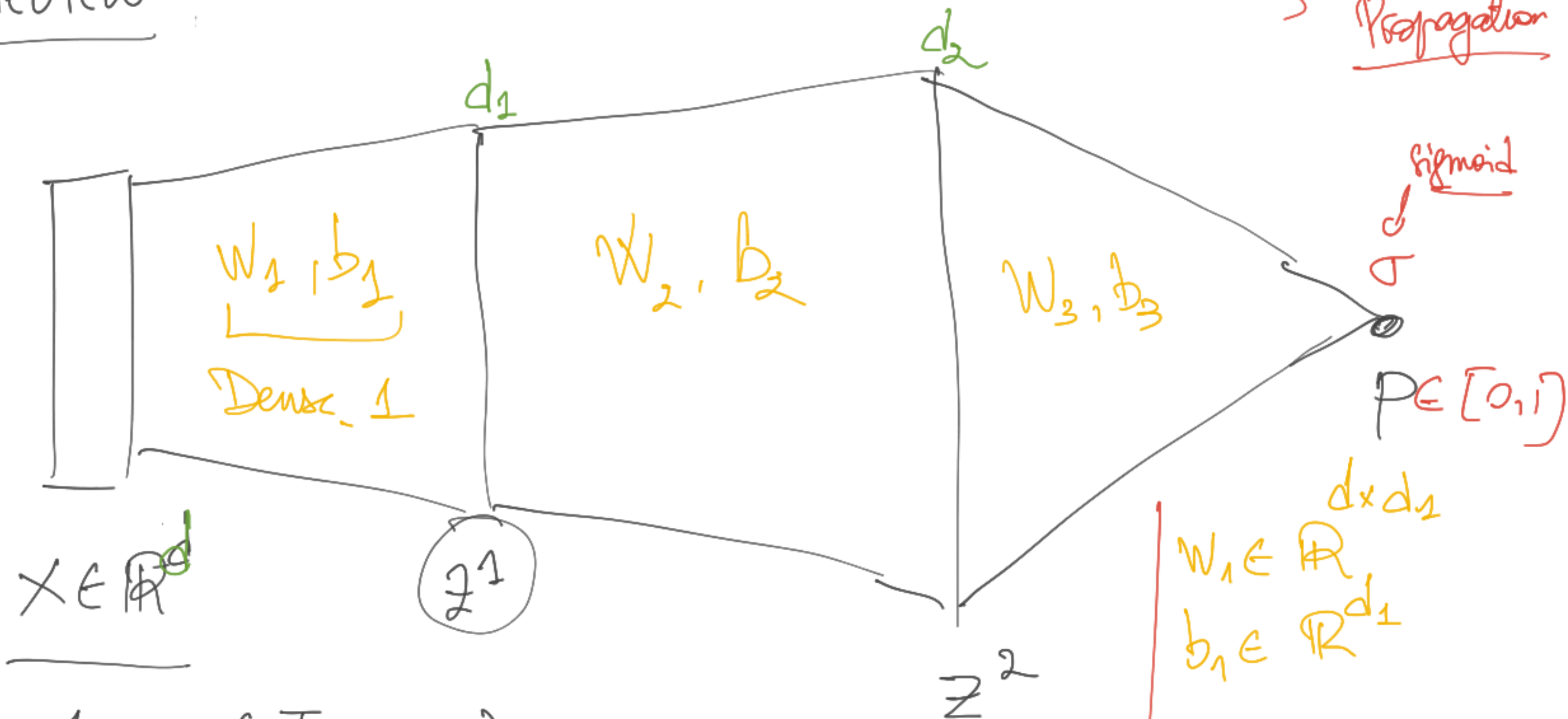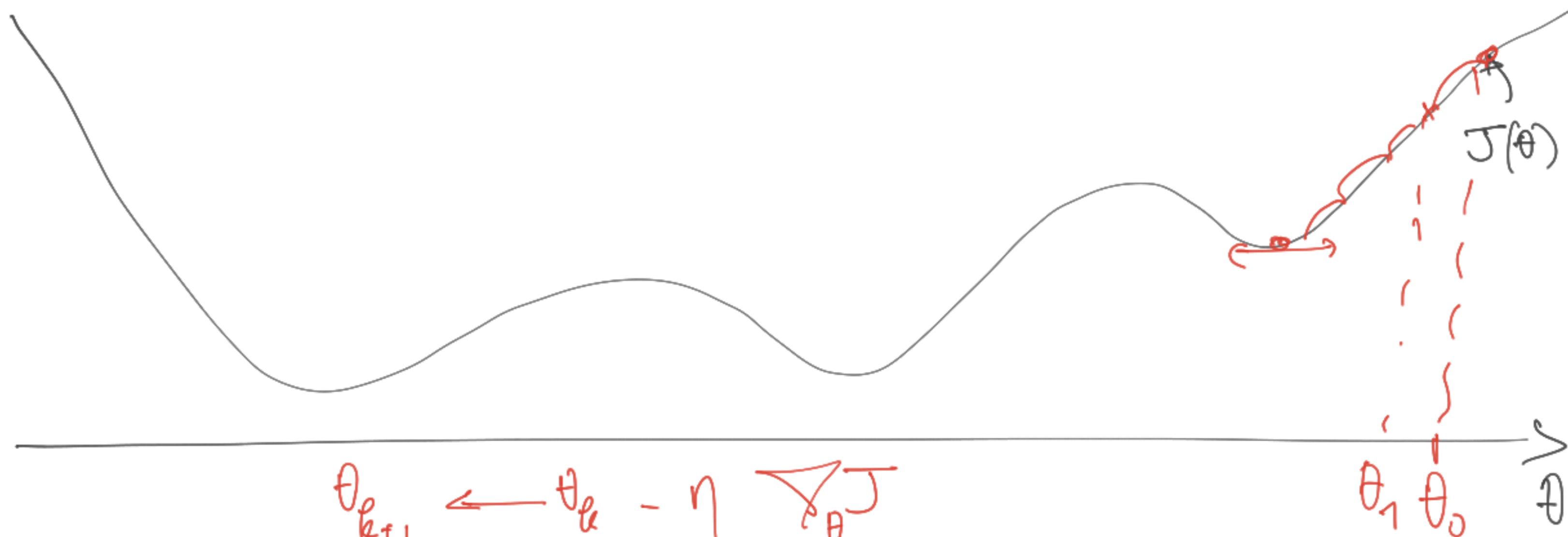$d \times d_1$

$W_1 \in \mathbb{R}^{d \times d_1}$

$b_1 \in \mathbb{R}^{d_1}$

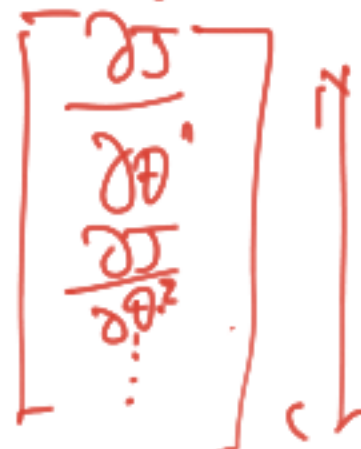$$z^1 = \sigma(W_1^T x + b_1)$$
$$z^2 = \sigma(W_2^T z^1 + b_2) \quad ; \quad P = \sigma(W_3^T z^2 + b_3)$$

Parameters : $\theta = \{ W_i, b_i \ ; \ 1 \leq i \leq 3 \} = \begin{bmatrix} \ \\ \ \end{bmatrix}$ all parameters



$J(\theta)$

$\theta_1 \ \theta_0$

$\theta$

$\theta_{k+1} \longleftarrow \theta_k - \eta \ \nabla_\theta J$

$\begin{bmatrix} \ \\ \ \end{bmatrix} = \begin{bmatrix} \ \\ \ \end{bmatrix} - \eta \underset{\sim 10^{-3}}{\downarrow} \begin{bmatrix} \frac{\partial J}{\partial \theta^1} \\ \frac{\partial J}{\partial \theta^2} \\ \vdots \end{bmatrix}$  # parameters

Loss function: $J(\theta) = -\dfrac{\log(L(\theta))}{N}$

↳ Binary Classification ⊕:

$$J(\theta) = \frac{-1}{N} \sum_{i=1}^{N} \left\{ y_i \log p_i + (1-y_i) \log(1-p_i) \right\}$$

Dataset: $\{\underset{d}{x_i}, y_i\}$, $1 \leq i \leq N$

$x_i \xrightarrow{\quad f_\theta \quad} p_i$

$p_i = \boxed{f_\theta}(x_i)$ ← NN

↳ **Multiclass classification pb** : (with $K$ categories)

Dataset : $\{x_i, y_i\}_{1 \leq i \leq N}$ ; $x_i \in \mathbb{R}^d$ ; $y_i \in \{1, \dots, K\}$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ ↳ $\hat{y}_i = [0 \dots 1 0 \dots 0]_K$

$$x_i \xrightarrow{\;f_\theta\;} \underline{\underline{P_i}} = \begin{bmatrix} P_i^1 \\ \vdots \\ P_i^K \end{bmatrix}$$

$$J(\theta) = \frac{1}{N} \sum_{i=1}^{N} \sum_{k=1}^{K} \hat{y}_i^k \log(P_i^k)$$

**Dataset :** Approach in Programming Session 4.

$doc_1$ : "Neural Networks are fun" $\Rightarrow$ [————————]$_V$
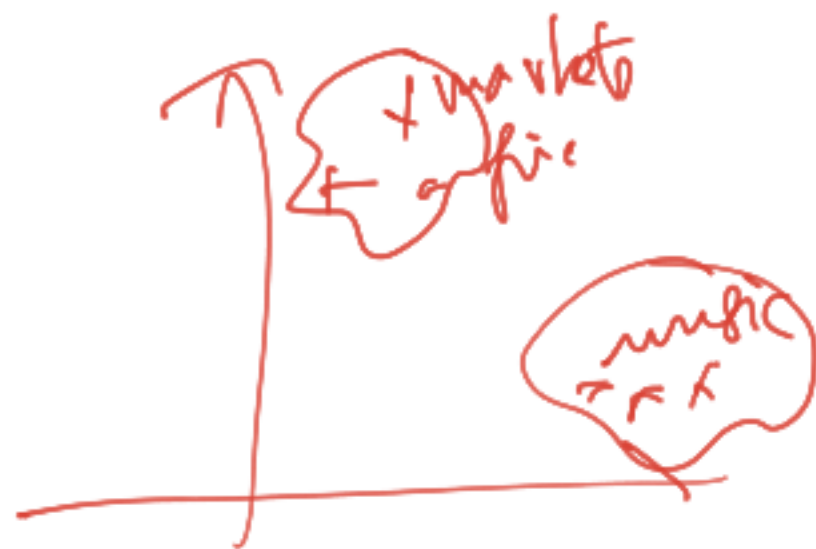
$doc_2$ : Attention models are awesome $\Rightarrow$ [————————]$_V$

$\vdots$

$doc_N$ : "This course is boring $\Rightarrow$ [————————]$_V$

---

data processed is of shape $(N, V)$

## Recap

### GloVe approach

$D \sim 50 - 300$

$$W = \begin{bmatrix} \rule[0.5ex]{1em}{0.4pt} & W_1 & \rule[0.5ex]{1em}{0.4pt} \\ & \vdots & \\ & \vdots & \\ \rule[0.5ex]{1em}{0.4pt} & W_V & \rule[0.5ex]{1em}{0.4pt} \end{bmatrix} \in \mathbb{R}^{V \times D}$$

}  D-dim representation of each word (as a center word)

(~ users)

$$\tilde{W} = \begin{bmatrix} \rule[0.5ex]{1em}{0.4pt} & \tilde{W}_1 & \rule[0.5ex]{1em}{0.4pt} \\ & \vdots & \\ \rule[0.5ex]{1em}{0.4pt} & \tilde{W}_V & \rule[0.5ex]{1em}{0.4pt} \end{bmatrix} \in \mathbb{R}^{V \times D}$$

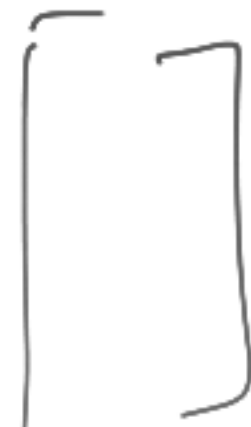}  D-dim representation of each word (as a context word)

(~ movies)

# Neural Networks as fun.

32          47          100          2



$e_{32}$       $e_{47}$       $e_{100}$       $e_2$

$\longrightarrow$ T

Embedding
Matrix
(GloVe)

$$e = \begin{bmatrix} - e_1 - \\ - e_2 - \\ \vdots \\ - e_v - \end{bmatrix} \begin{matrix} 1 \\ \\ \vdots \\ v \end{matrix}$$

$\in \mathbb{R}^{V \times D}$

$doc_1 \longrightarrow [\, e_{32},\ e_{47},\ e_{100},\ e_2 \,] \in \mathbb{R}^{\underline{T \times D}}$

$\begin{bmatrix} \ \end{bmatrix}_D \quad \begin{bmatrix} \ \end{bmatrix}_D \quad \begin{bmatrix} \ \end{bmatrix}_D \quad \begin{bmatrix} \ \end{bmatrix}_D$

$$doc_1 \longrightarrow (T, D)$$

$$doc_2 \longrightarrow (T, D)$$

$$\vdots$$

$$doc_N \longrightarrow (T, D)$$

$$\Bigg\} \quad \text{Dataset processed} \quad (N, T, D)$$