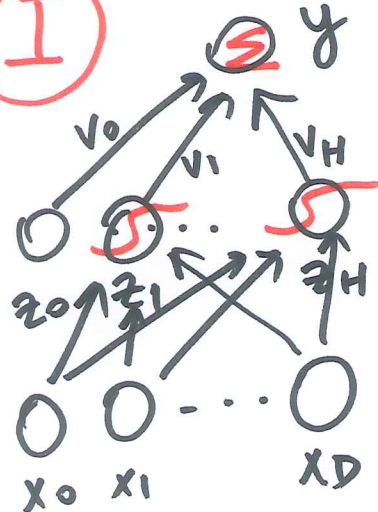


# MULTILAYER PERCEPTRONS

I



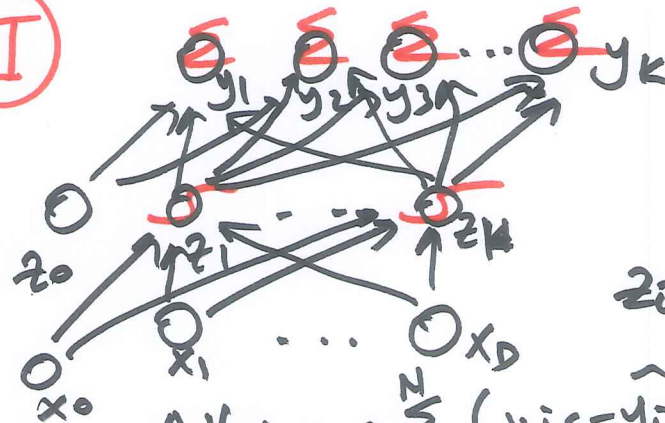
$$\hat{y}_i = \underline{v}^T \cdot z_i$$

$$z_{ih} = \text{sigmoid}(\underline{w}_h^T \cdot x_i)$$

$$\underline{\Delta v}_h = \eta \sum_{i=1}^N (y_i - \hat{y}_i) \cdot z_{ih}$$

$$\underline{\Delta w}_{hd} = \eta \sum_{i=1}^N (y_i - \hat{y}_i) \cdot v_h \cdot z_{ih} (1 - z_{ih}) \cdot x_{id}$$

II



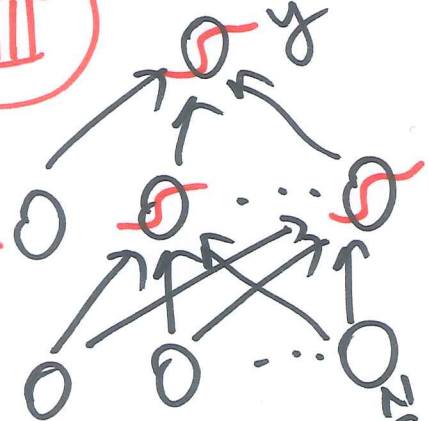
$$\hat{y}_{ic} = \underline{v}_c^T \cdot z_i$$

$$z_{ih} = \text{sigmoid}(\underline{w}_h^T \cdot x_i)$$

$$\underline{\Delta v}_{ch} = \eta \sum_{i=1}^N (y_{ic} - \hat{y}_{ic}) \cdot z_{ih}$$

$$\underline{\Delta w}_{hd} = \eta \sum_{i=1}^N \left[ \sum_{c=1}^K (y_{ic} - \hat{y}_{ic}) v_{ch} \right] \cdot z_{ih} (1 - z_{ih}) \cdot x_{id}$$

III



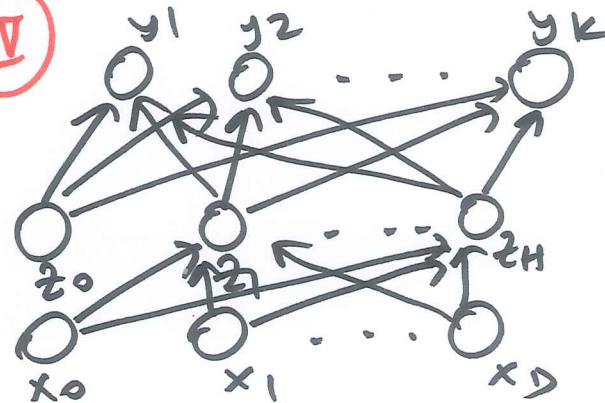
$$\hat{y}_i = \text{sigmoid}(\underline{v}^T \cdot z_i)$$

$$z_{ih} = \text{sigmoid}(\underline{w}_h^T \cdot x_i)$$

$$\underline{\Delta v}_h = \eta \sum_{i=1}^N (y_i - \hat{y}_i) \cdot z_{ih}$$

$$\underline{\Delta w}_{hd} = \eta \sum_{i=1}^N (y_i - \hat{y}_i) v_h \cdot z_{ih} (1 - z_{ih}) \cdot x_{id}$$

IV



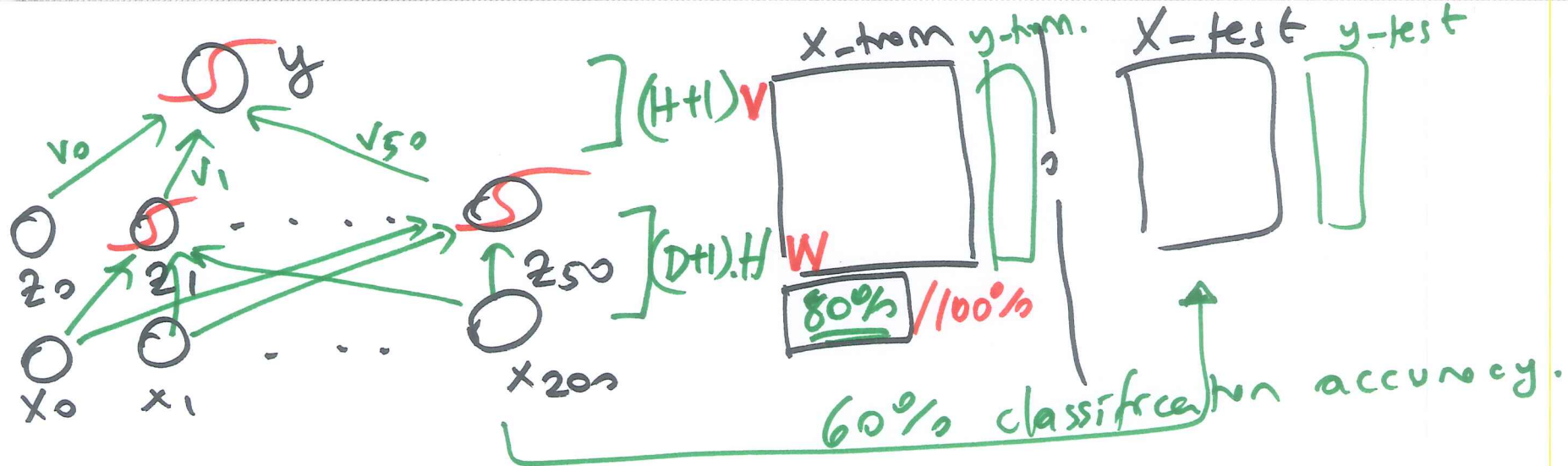
$$\hat{y}_{ic} = \text{softmax}(\underline{v}_c^T \cdot z_i)$$

$$\hat{y}_{ic} = \frac{\exp(\underline{v}_c^T \cdot z_i)}{\sum_{d=1}^K \exp(\underline{v}_d^T \cdot z_i)}$$

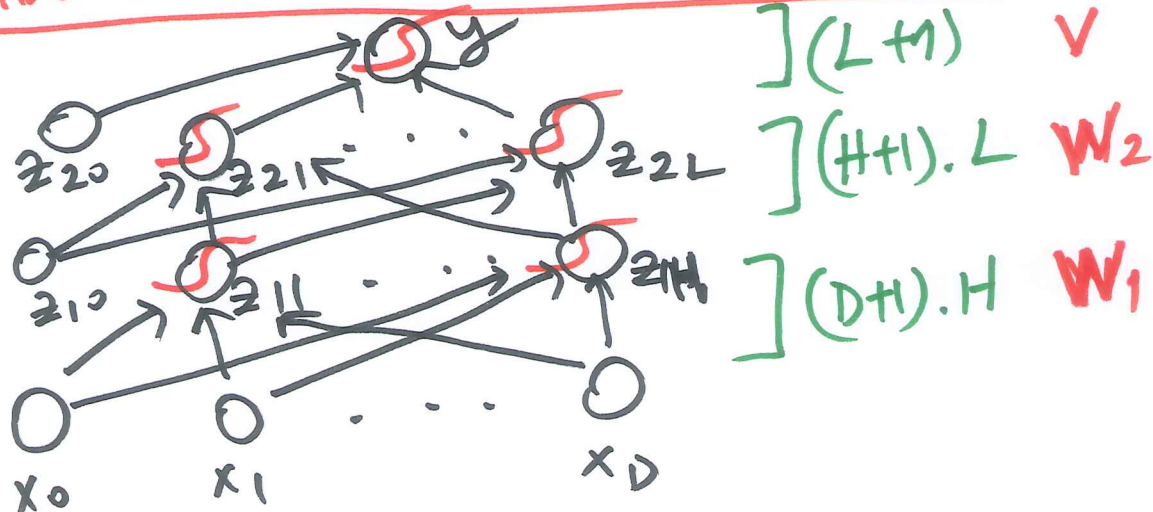
$$\underline{\Delta v}_{ch} = \eta \sum_{i=1}^N (y_{ic} - \hat{y}_{ic}) \cdot z_{ih}$$

$$\underline{\Delta w}_{hd} = \eta \sum_{i=1}^N \left[ \sum_{c=1}^K (y_{ic} - \hat{y}_{ic}) \cdot v_{ch} \right] \cdot z_{ih} (1 - z_{ih}) \cdot x_{id}$$

①



## MULTIPLE HIDDEN LAYERS



$$y = \text{sigmoid}(v^T \cdot z_2)$$

$$z_{2L} = \text{sigmoid}(w_{2L}^T \cdot z_1)$$

$$z_{1h} = \text{sigmoid}(w_{1h}^T \cdot x)$$

## Training Procedures

Momentum:

$$\Delta w_h^{(+)} = -\eta \frac{\partial \text{Error}^{(+)}}{\partial w_h} + \underbrace{\alpha \Delta w_h^{(+ - 1)}}_{\text{momentum term}}$$

$$0.5 < \alpha < 1$$

Adaptive Learning Rate:

$\eta \Rightarrow$  should be getting smaller in the last iterations.

Early stopping:

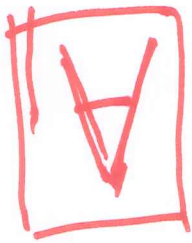
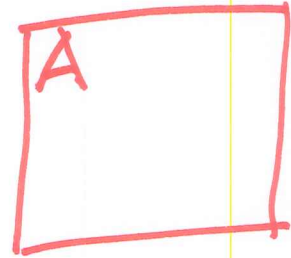
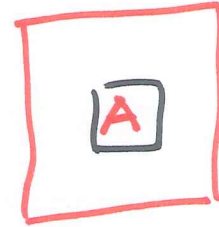
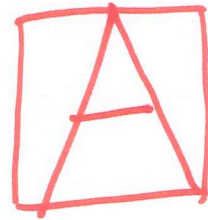
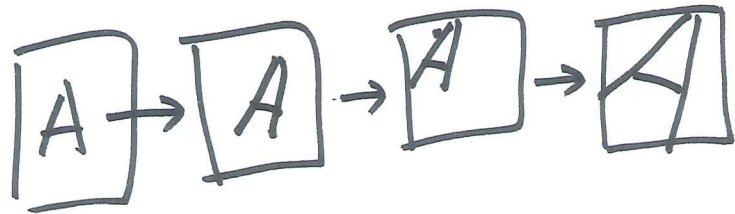




# HINTS

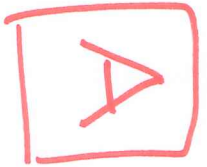
prior knowledge (hints)

A Z B II B



scaled

translated  
scaled



⋮

## ① virtual examples:

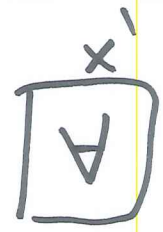
- randomly perturb your training data
- add these perturbed data points to your training data set

## ② preprocessing: Just train your algorithm with preprocessed images.

## ③ special network structure: (weight sharing).

## ④ modifying error function:

$$\text{Error}' = \text{Error} + [g(x|\theta) - g(x'|\theta)] \underline{g(x|\theta)}$$



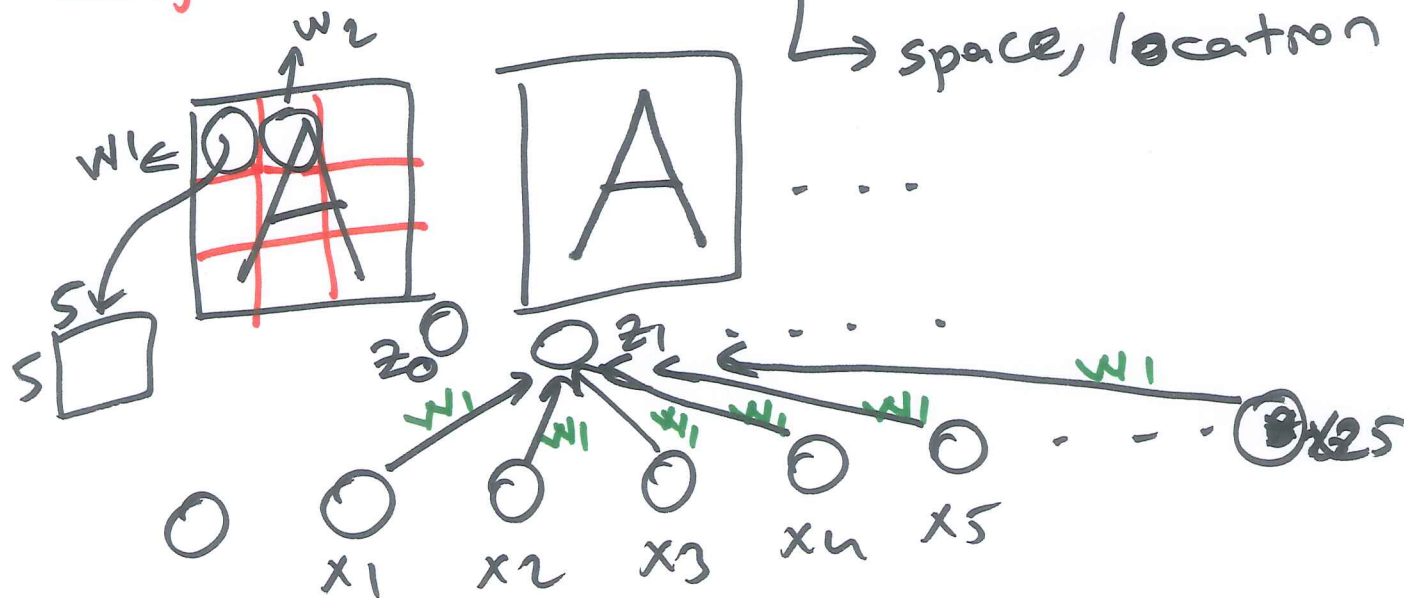
$$\underline{g(x'|\theta)}$$

④

weight sharing:

spatial dependency

↳ space, location



temporal dependency  
↳ time.

