

#3) i) Space (Storage) complexity :

of adjustable weights from input to hidden layer is : $I \cdot J$

of adjustable weights from hidden to output layer is : $J \cdot K$

Assuming thresholds are all at zero, total # weights = $I \cdot J + J \cdot K$

During learning, however, one also has to store the gradients \Rightarrow total storage for weights + their gradients is : $2(IJ + JK)$.

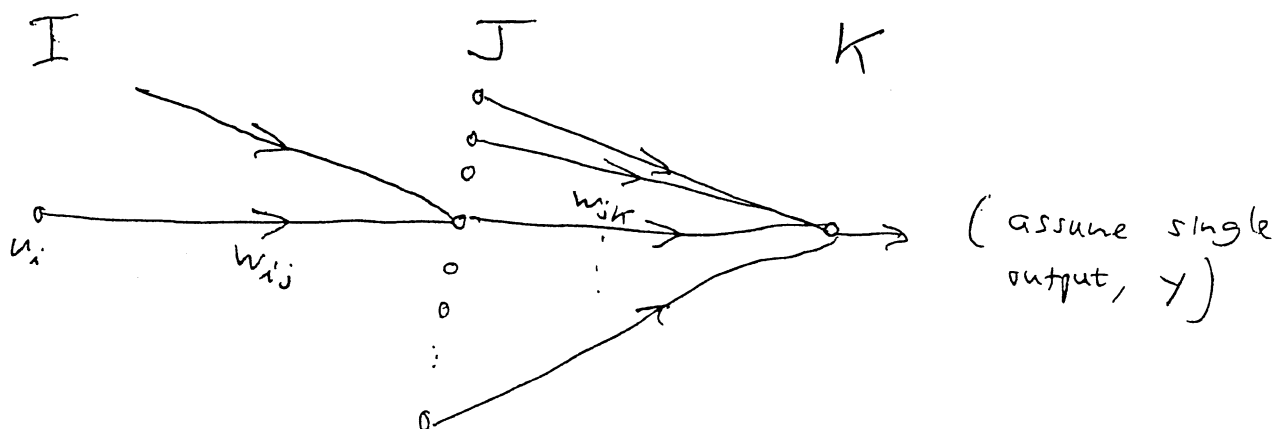
We also need to store the training patterns \Rightarrow

$(I+1)T$
↓ ↓ ↘
input output # of training
dim. examples

\Rightarrow space complexity = $2(IJ + JK) + (I+1)T$

ii) Computational complexity of BP :

Refer to the BP derivation from lecture :



$$V_j = \sum_{i=1}^I w_{ij} u_i, \quad j=1, \dots, J \Rightarrow O(IJ) \text{ operations, per sample}$$

$$y = g\left(\sum_{j=1}^J w_{j0} g(V_j)\right) \Rightarrow O(J) \text{ operations, per sample}$$

Forward operations

(mults. & adds.)

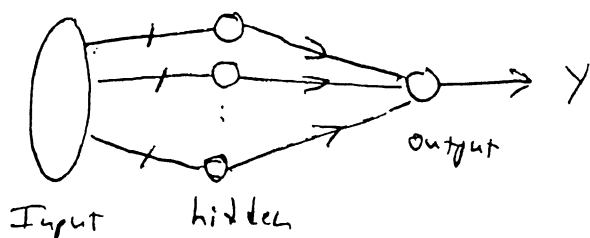
$$\delta_j = g'(V_j) \delta_0 w_{j0}, \quad j=1, \dots, J \Rightarrow O(J) \text{ operations}$$

$$\Delta w_{ij} \propto \delta_j \cdot u_i, \quad i=1, \dots, I, j=1, \dots, J \Rightarrow O(IJ) \text{ operations}$$

$$\Delta w_{j0} \propto -\frac{\partial E}{\partial V_0} \cdot \frac{\partial V_0}{\partial w_{j0}}, \quad j=1, \dots, J \Rightarrow O(J) \text{ operations}$$

Overall complexity is $O(I \cdot J \cdot T)$ operations

#4) Consider the MLP:



Suppose each neuron uses an antisymmetric activation function, i.e. $f(-x) = -f(x)$ (sigmoid, tanh, or ± 1).

The output can be written

$$Y = f\left(\sum_{j=1}^{N_h} \underbrace{w_{j0}}_{\text{output}} f\left(\sum_{i=1}^d x_i w_{ij} + w_{j0}\right) + w_0\right)$$

Spec. we flip the sign on all the weights.

Then, neuron j 's output becomes

$$f\left(\sum_{i=1}^d x_i (-w_{ij}) + (-w_{j0})\right) = -f\left(\sum_{i=1}^d x_i w_{ij} + w_{j0}\right)$$

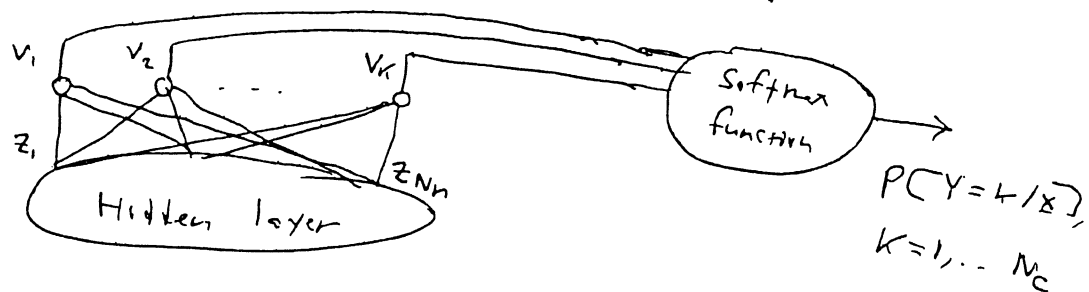
The overall output becomes:

$$\begin{aligned} & f\left(\sum_{j=1}^{N_h} (-w_{j0}) (-f\left(\sum_{i=1}^d x_i w_{ij} + w_{j0}\right)) - w_0\right) \\ &= f\left(\sum_{j=1}^{N_h} w_{j0} f\left(\sum_{i=1}^d x_i w_{ij} + w_{j0}\right) - w_0\right). \end{aligned}$$

The output stays the same if $w_0 = 0$ -- otherwise the operation is in fact changed.

This property does NOT extend to flipping the sign on the input.

#5) Let's just consider the signals from the hidden layer to the softmax function:



Let's denote the outputs of the hidden layer by:

$$z_j, j=1, \dots, N_h.$$

$$\text{Then, } v_k = \sum_{j=1}^{N_h} z_j w_{jk} \quad \text{and } P(Y=k/X) = \frac{e^{v_k}}{\sum_{k'=1}^{N_c} e^{v_{k'}}}$$

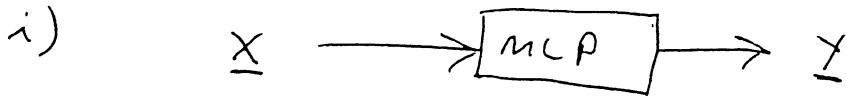
As derived in lecture, the cross entropy objective function can be written as:

$$E = - \sum_{c=1}^{N_c} \sum_{\underline{x}: c(\underline{x})=c} \log P(Y=c/X)$$

$$\text{Now, } \Delta w_{jk} \propto \frac{\partial E}{\partial v_k} \cdot \underbrace{\frac{\partial v_k}{\partial w_{jk}}}_{z_j}$$

$$\begin{aligned} \frac{\partial E}{\partial v_k} &= - \sum_{\underline{x}: c(\underline{x})=k} \frac{\partial}{\partial v_k} (\log P(Y=k/X)) - \sum_{\substack{c=1, \\ c \neq k}}^{N_c} \sum_{\underline{x}: c(\underline{x})=c} \frac{\partial}{\partial v_k} (\log P(Y=c/X)), \\ \text{where:} \quad \frac{\partial}{\partial v_k} (\log P(Y=k/X)) &= \frac{\partial}{\partial v_k} \left(v_k - \log \left(\sum_{k'=1}^{N_c} e^{v_{k'}} \right) \right) \\ &= 1 - \frac{e^{v_k}}{\sum_{k'=1}^{N_c} e^{v_{k'}}} = 1 - P(Y=k/X) \end{aligned}$$

#6) Neural Network Inversion



$\underline{y} = f(\underline{x})$, with $f(\cdot)$ defined by the MLP. Ideally, if we were given a vector \underline{y} , we could find

\underline{x} via $\underline{x} = f^{-1}(\underline{y})$. Unfortunately,

- i) there are ~~many~~ ^{may be} possible \underline{x} that could lead to same $\underline{y} \Rightarrow$ non-unique inverse.
- ii) there may be no \underline{x} leading to a particular \underline{y} .
- iii) In general, there is no analytical form for $f^{-1}(\cdot)$.

So, how to find an approximate inverse?

Ans: Use backpropagation!

Let $E = \|\underline{y} - f(\underline{x})\|^2$ + choose \underline{x} to minimize

E . Algorithm:

$t=0; \underline{x}^t = \underline{x}_0$ (Init) \longleftarrow

Note: strong dependence on initialization...

loop

$$\underline{x}_{t+1} = \underline{x}_t - \eta \nabla_{\underline{x}} E(\underline{x}) \Big|_{\underline{x} = \underline{x}_t}$$

$t \leftarrow t+1$

end loop

- Two applications:
- 1) signal/image reconstruction from noisy, blurred output image.
 - 2) Boundary-finding in classification -- given an MLP classifier, find \underline{x} that cause output to be $\approx 1/2$ (points on the boundary ...).