

Lecture 8:

Feed-forward neural networks (FFNNs)

CS 4740 (and crosslists): Introduction to Natural Language Processing

Slides developed by:

Magd Bayoumi, Claire Cardie, Tanya Goyal, Dan Jurafsky, Lillian Lee, James Martin, Marten van Schijndel

Announcements

- **HW0 and HW1 milestone graded on Gradescope
HW1 due on Gradescope on Friday by 11:59pm**

So far this semester...computational models + NLP tasks

- N-gram language models
 - Predict the probability of a sequence of tokens
 - Predict the next token in a sequence
- HMMs for sequence tagging tasks
 - E.g., POS tagging, NER
 - MEMMs (briefly covered)
 - $P(t_i | t_{i-1}) \rightarrow P(t_i | t_{i-1}, o_{1:N})$
- Logistic regression for text classification (binary)
 - E.g., document/sentence-level sentiment classification
 - Convert input into a set of *feature-value* pairs
 - Learns *weight* to apply to each feature-value pair

Last class: Representing word types

- How did we represent words for each of the NLP models we've covered?

Recall from the last class

- Words as vectors

Vectors of words with similar meanings/uses should be close in the vector space

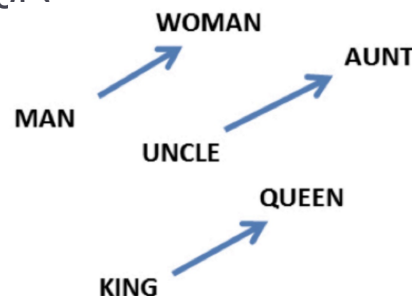
Co-occurrence vector that encodes *context* using raw counts

| | | | | | | | | |
|-------------|----------|-----|----------|------|--------|-----|-------|-----|
| | aardvark | ... | computer | data | result | pie | sugar | ... |
| digital | 0 | ... | 1670 | 1683 | 85 | 5 | 4 | ... |
| information | 0 | ... | 3325 | 3982 | 378 | 5 | 13 | ... |

Improve using *tf-idf* representation


Word embeddings – (short) dense word vectors

- Learn via the *skip-gram* model (Word2Vec)
- Each vector element, $0 < v_i < 1$
- Nice properties naturally emerge (but so do *biases* from the data)



Immediately improved performance on virtually all NLP tasks!!!

Today

- Feed-forward neural networks
 - As a mechanism for *representation learning*  **for word embeddings**
 - As another approach to classification problems for text/language

The problem: generic word vectors might not be optimal

For a *specific* task, such as NE tagging, we could learn better representations.

- Example: it'd be nice if the vectors for LOC words were separate from the vectors for PER words.

The problem: generic word vectors might not be optimal

Given: word2vec rep for word w ,

$$\mathbf{x} = [x_1, \dots, x_i, \dots, x_{n_0}]$$



new vector rep
 $\mathbf{h} = [h_1, \dots, h_{n_1}]$



Predict: set of n_2 task-specific outputs,

$$\mathbf{y} = [y_1 = P(\text{ORG}|w), y_2 = P(\text{PER}|w), \dots, y_{n_2} = P(\text{LOC}|w)]$$

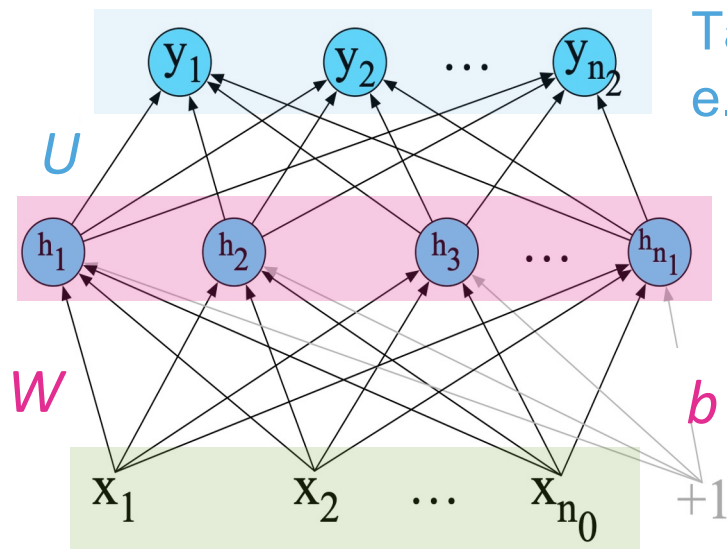
FFNNs will acquire the ability to perform this task by learning to derive better representations of the word vectors.

The solution we're building up to today: feed-forward neural networks (FFNNs) for classification

- Network of small simple(!) computing units
 - Each takes a vector of input values and produces a single output value
- Computation for classification proceeds from one layer of units to the next
- **Deep learning** when the network has many layers
- Logistic regression is a really basic FFNN
- Unlike LR, **no feature engineering** is needed for FFNNs!!!
 - Take dense word embeddings as input
 - Induce the necessary features as part of learning to classify
- Downside: need a large amount of training data



The solution we're building up to today: feed-forward neural networks (FFNNs)



Task-specific outputs,
e.g., $y_1 = P(\text{ORG}|w)$

new vector rep
 $\mathbf{h} = [h_1, \dots, h_{n_1}]$

for word w , word2vec input
 $\mathbf{x} = [x_1, \dots, x_i, \dots, x_{n_0}]$

Intuition: FFNNs learn how to better represent \mathbf{x} (i.e, as \mathbf{h}) so as to be able to accurately predict \mathbf{y} .

The problem again: generic word vectors might not be optimal for a specific task

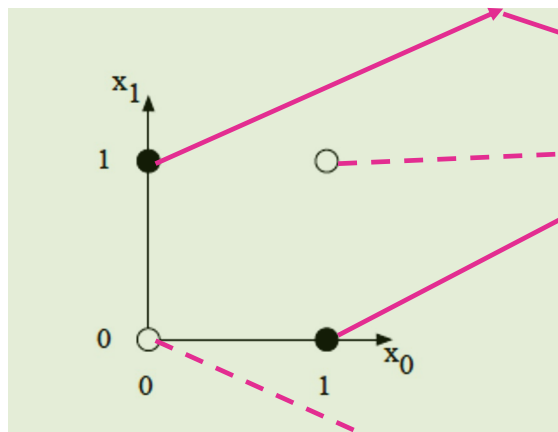
Recall: for a *specific* task, such as NE tagging, we could learn better representations.

- Example: it'd be nice if the vectors for ORG words were separate from the vectors for PER words.

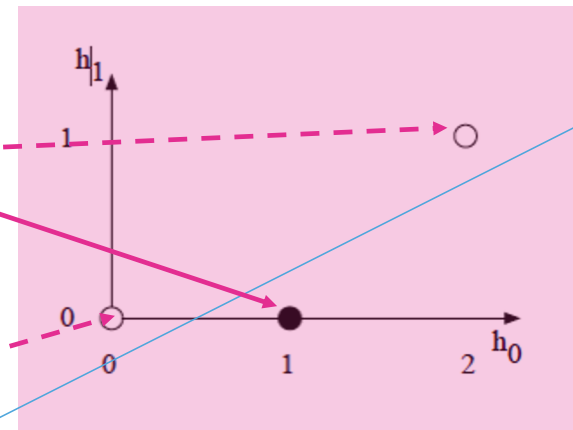
Changing the original x 's into new h 's where items with different labels are separable

Suppose that \bullet = PER, and \circ = LOC, and our $w2v$ vectors looked like this:

There's a transform whereby the PERs can be separated from the LOCs:

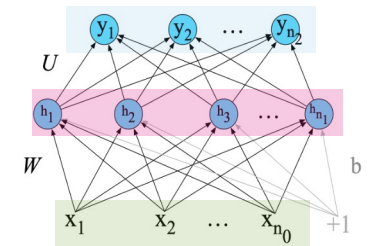
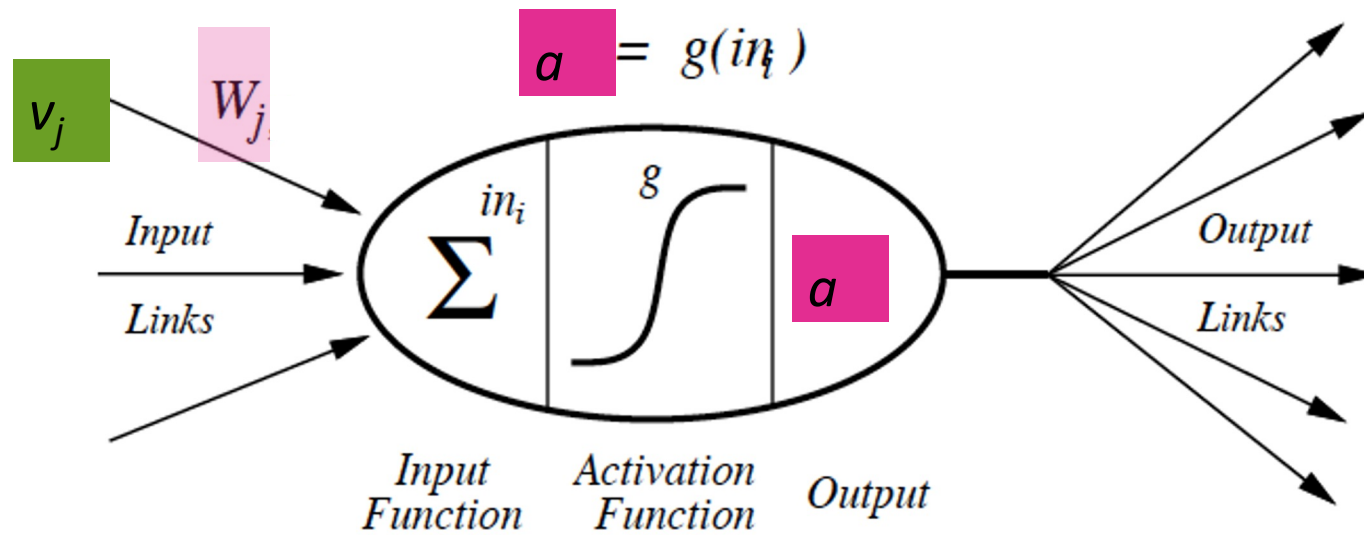


a) The original x space



b) The new h space

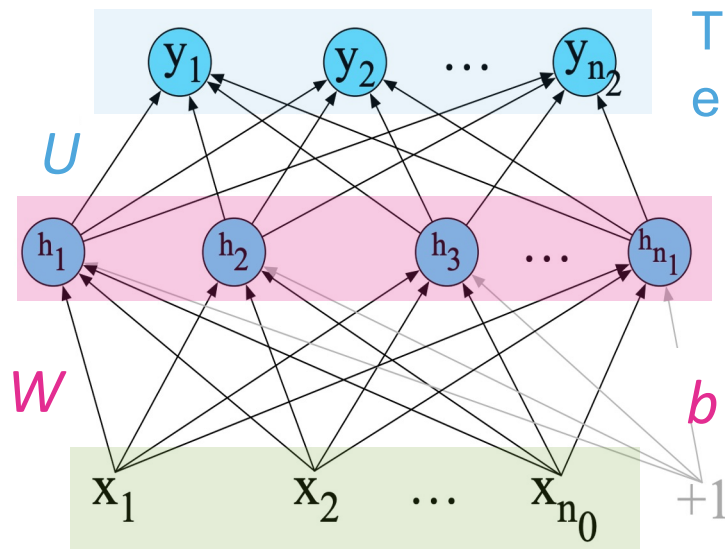
Our transformations will use
“computational neurons”/units.



Or, in other words, the scalar $a := g(\sum_j w_j v_j)$.

Or, if we let $\mathbf{w} = [\dots, w_j, \dots]$ and $\mathbf{v} = [\dots, v_j, \dots]^T$, $a := g(\mathbf{w}\mathbf{v})$

Feed-forward neural networks (FFNNs)



Task-specific outputs,
e.g., $y_1 = P(\text{ORG}|w)$

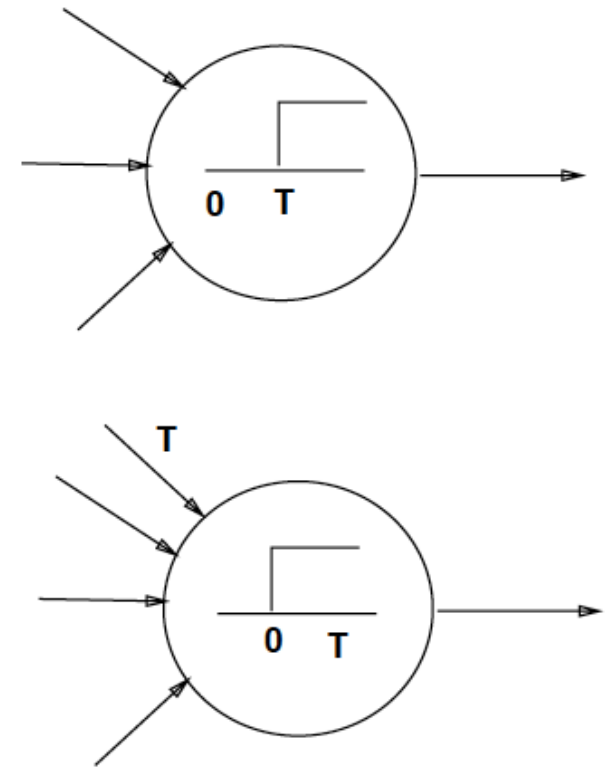
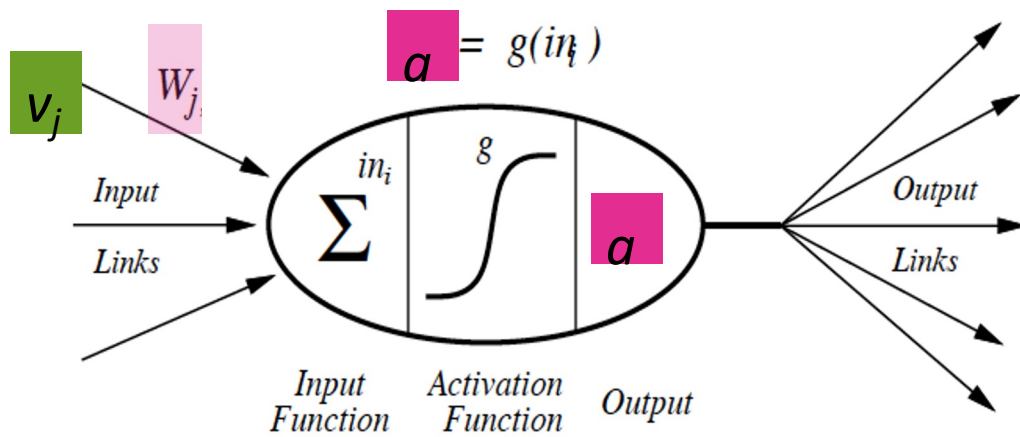
During training, all of the weights will be learned.

for word w , word2vec input
 $\mathbf{x} = [x_1, \dots, x_i, \dots, x_{n_0}]$

What is this b term and the extra input value that is always a 1?

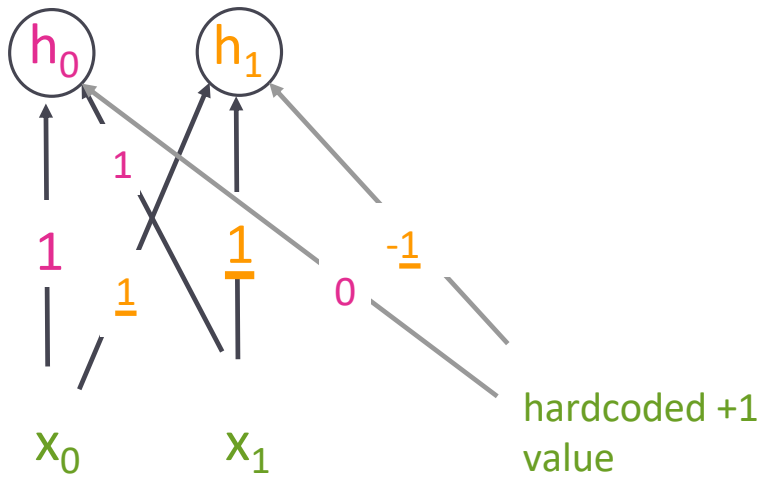
The *bias* term

Allows learning of the thresholds as though they were weights.

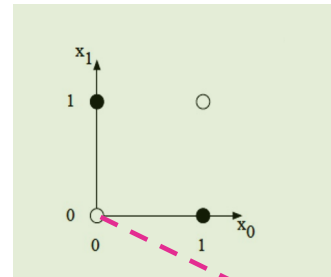


Our LOC/PER example

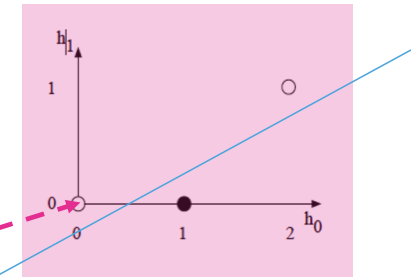
g = the ReLU activation function,
 $g(x) = \max(x, 0)$



Trick: This hardcoded +1 input allows us to treat the bias scalar b as just another weight.



a) The original x space



b) The new h space

Computation for $[0, 0]$, \circ = LOC:

$$h_0 = g(1(0) + 1(0) + 0(+1))$$

$$= g(0) = 0$$

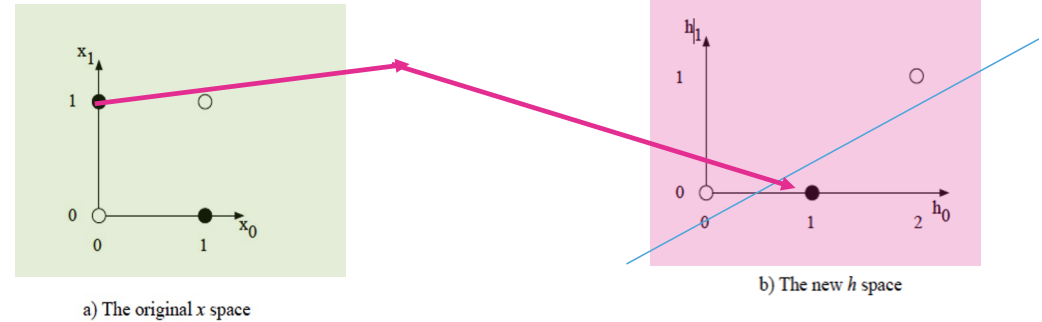
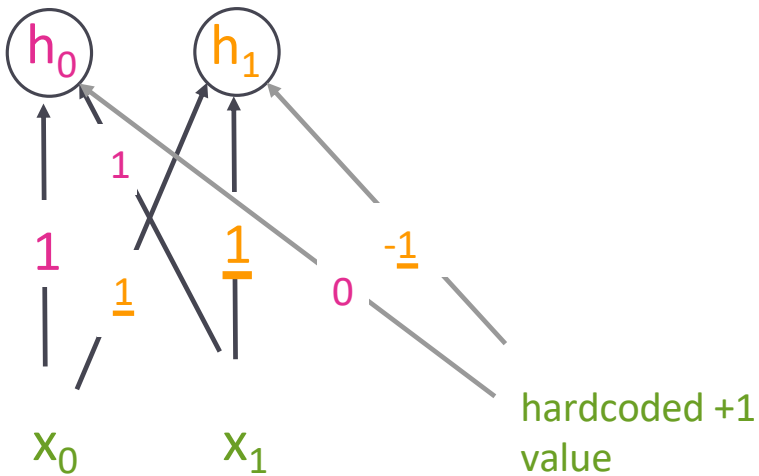
$$h_1 = g(\underline{1}(0) + \underline{1}(0) + -1(+1))$$

$$= g(-1)$$

$$= 0$$

Our LOC/PER example

g = the ReLU function,
 $g(x) = \max(x, 0)$



Computation for $[0, 1]$, $\square = \text{PER}$

$$h_0 = g(1(0) + 1(1) + 0(+1))$$

$$= g(1) = 1$$

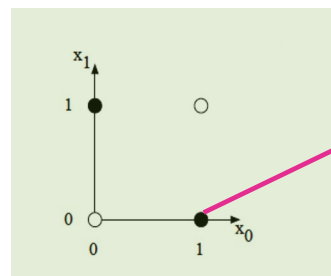
$$h_1 = g(\underline{1}(0) + \underline{1}(1) + -1(+1))$$

$$= g(0)$$

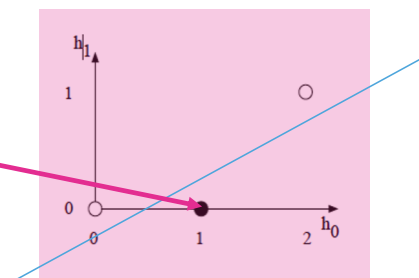
$$= 0$$

Switch to matrix computations!

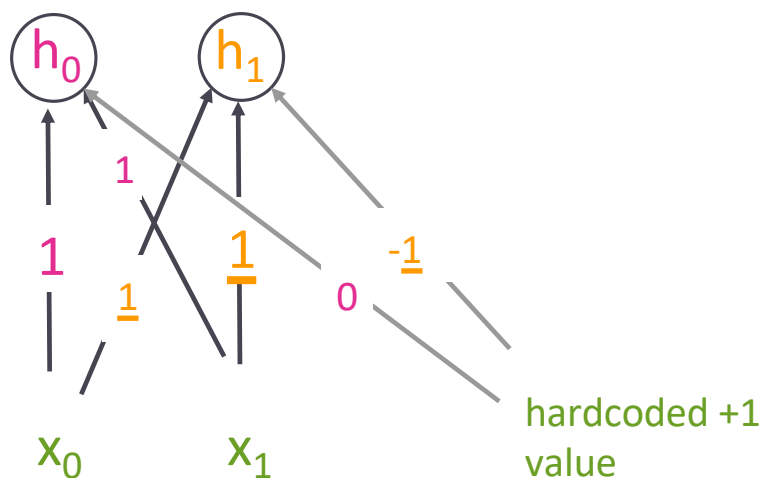
g = the ReLU function,
 $g(x) = \max(x, 0)$



a) The original x space



b) The new h space



Computation for $[1, 0]^T$,  = PER

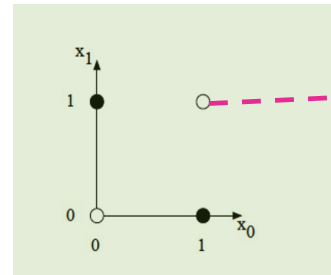
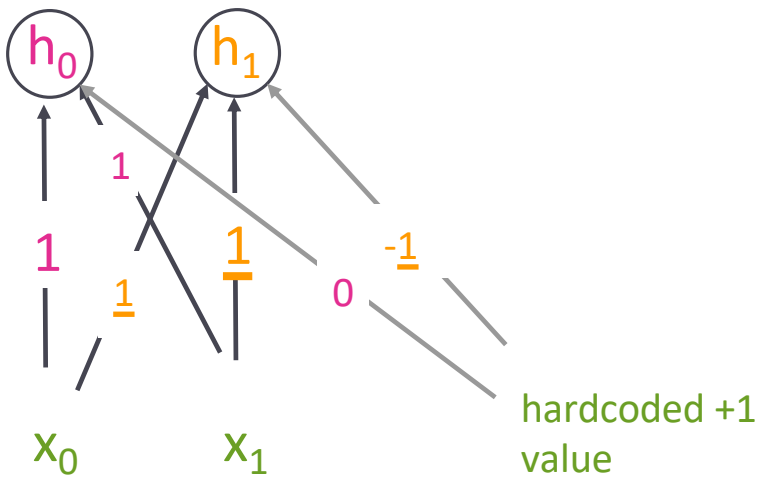
$$h_0 = g([1, 1, 0] [1, 0, +1]^T) \\ = g(1) = 1$$

$$h_1 = g([\underline{1}, \underline{1}, \underline{-1}] [1, 0, +1]^T) \\ = g(0) \\ = 0$$

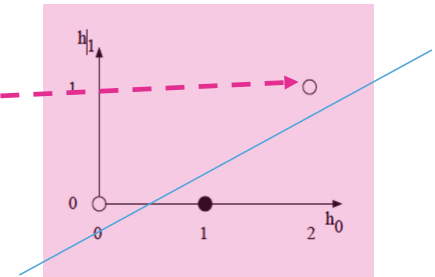
T = transpose: our word vectors are column vectors

Switch to matrix computations!

g = the ReLU function,
 $g(x) = \max(x, 0)$



a) The original x space



b) The new h space

Computation for $[1, 1]^T$, $\circ = \text{LOC}$

$$h_0 = g(1(1) + 1(1) + 0(+1)) = g(2) = 2$$

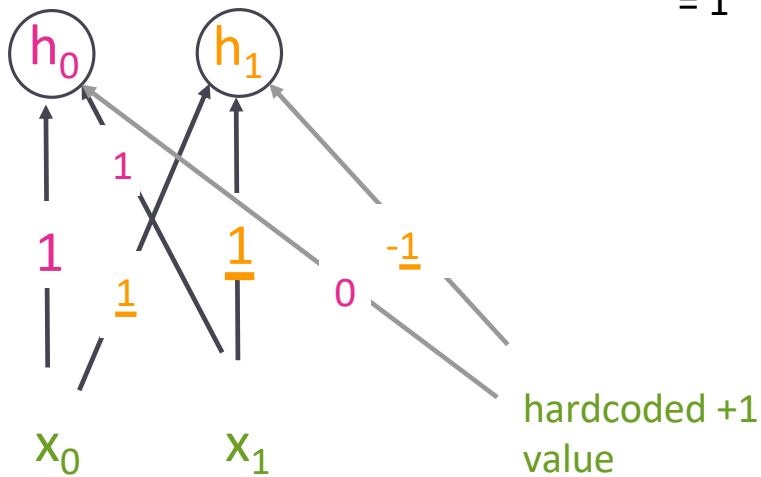
$$h_1 = g(\underline{1}(1) + \underline{1}(1) + -1(+1)) = g(1) = 1$$

Matrix version (use Python packages, not for-loops)

$$g \left(\begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & -1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ +1 \end{bmatrix} \right) = g \left(\begin{bmatrix} 2 \\ 1 \end{bmatrix} \right) = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

But, various computational packages want inputs as row vectors.
Be flexible!

g = the ReLU function,
 $g(x) = \max(x, 0)$



Computation for $[1, 1]^T$, $\circ = \text{LOC}$

$$h_0 = g(1(1) + 1(1) + 0(+1)) = g(2) = 2$$

$$h_1 = g(\underline{1}(1) + \underline{1}(1) + -1(+1)) = g(1) = 1$$

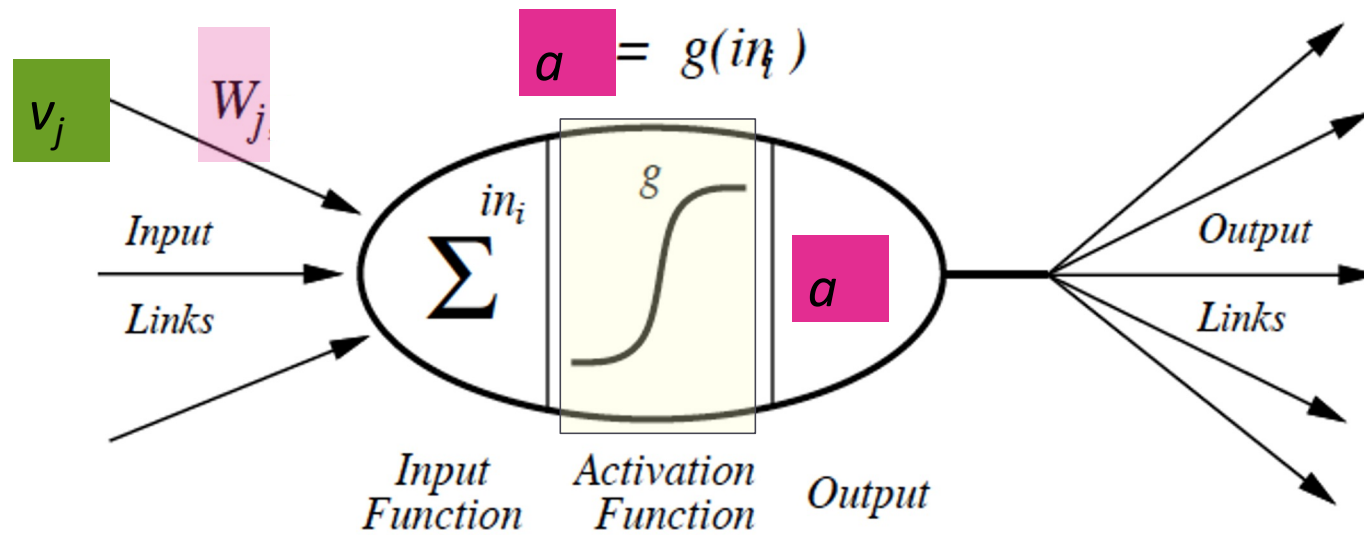
$$g \left(\begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & -1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ +1 \end{bmatrix} \right) = g \left(\begin{bmatrix} 2 \\ 1 \end{bmatrix} \right) = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

Is this equivalent?

row vector not column!

$$g \left(\begin{bmatrix} 1 & 1 & +1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 0 & -1 \end{bmatrix} \right) = g \left(\begin{bmatrix} 2 & 1 \end{bmatrix} \right) = \begin{bmatrix} 2 \\ 1 \end{bmatrix}^T$$

Alternate activation functions $g()$



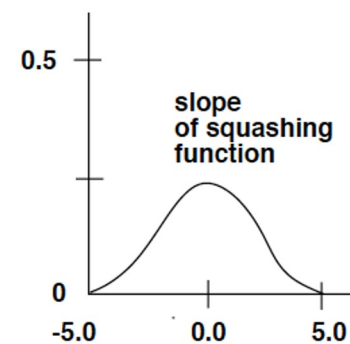
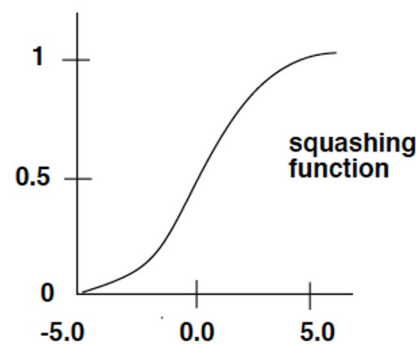
Or, in other words, the scalar $a := g(\sum_j w_j v_j)$.

Or, if we let $\mathbf{w} = [\dots, w_j, \dots]$ and $\mathbf{v} = [\dots, v_j, \dots]^T$, $a := g(\mathbf{w}\mathbf{v})$

Sigmoid activation function

- Maps the output into the range [0, 1]
 - useful in squashing outliers toward 0 or 1
- Smoothly **differentiable**: the training algorithm requires a differentiable activation function

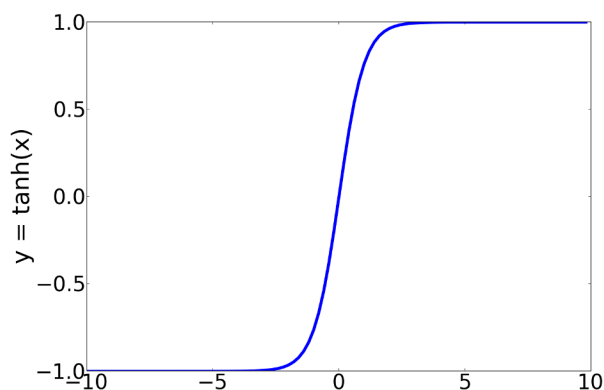
$$y = \sigma(z) = \frac{1}{1 + e^{-z}}$$



Other activation functions

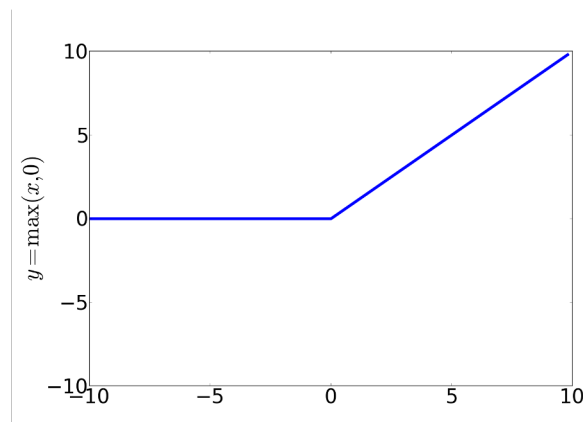
- ***tanh*** a variant of the sigmoid that ranges from -1 to +1

$$y = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

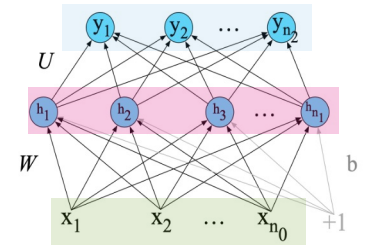


- ***ReLU*** rectified linear unit: equals x when x is positive, and 0 otherwise

$$y = \max(x, 0)$$



Feedforward network



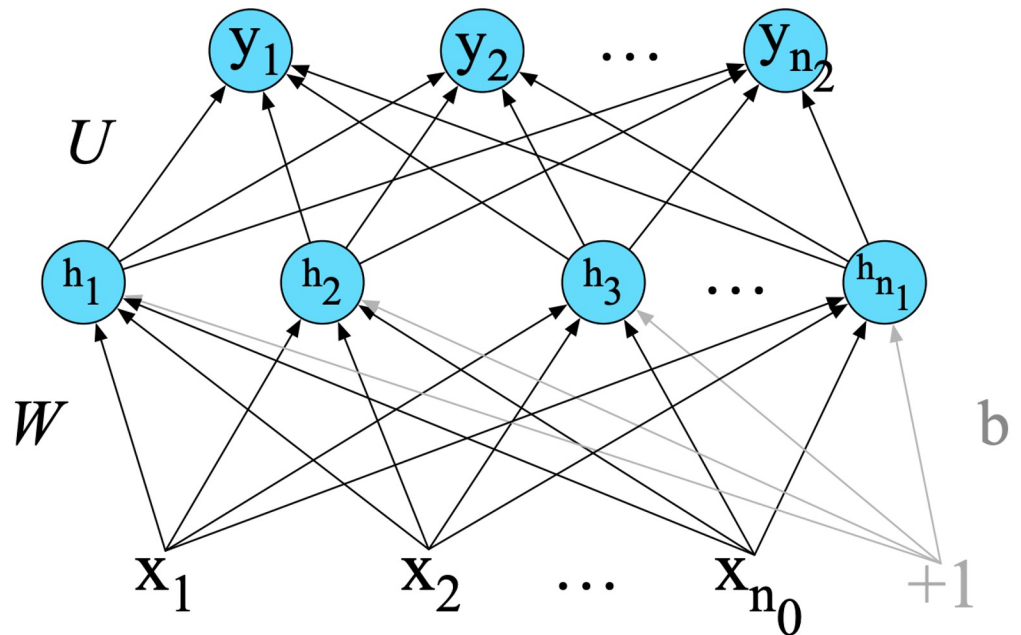
- Computation proceeds iteratively from one layer of units to the next.
- A feedforward network is a **multilayer** feedforward network in which the units are connected with **no cycles**.
I.e., the outputs from units in each layer are passed to units in the next higher layer, and no outputs are passed back to lower layers.

Feedforward network

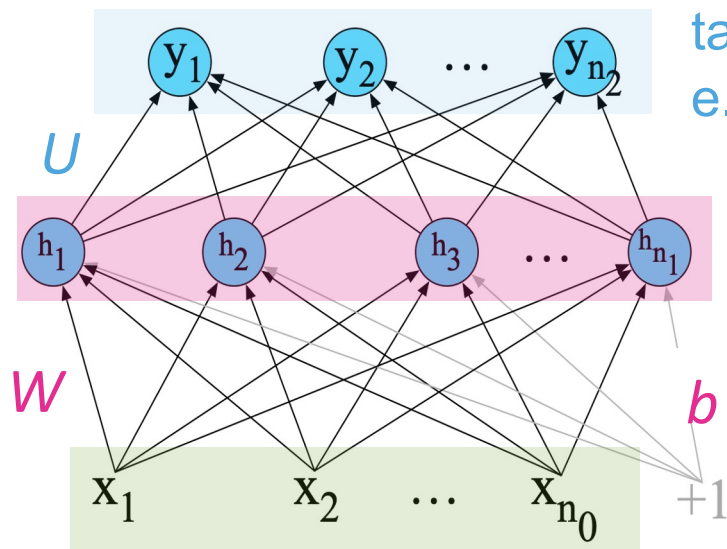
- Sometimes called **multi-layer perceptrons (MLPs)**; however, units in modern multilayer networks aren't perceptrons (linear)
They are units with non-linear activation functions (e.g. tanh).

Fully-connected FFNN

Each unit in each layer takes as input the outputs from all the units in the previous layer, and there is a link between every pair of units from two adjacent layers.



Back to the original picture! From \mathbf{x} to \mathbf{h} .



task specific outputs,
e.g., $y_1 = P(\text{ORG}|w)$

new vector rep
 $\mathbf{h} = [h_1, \dots, h_{n_1}]^T$

for word w , word2vec input
 $\mathbf{x} = [x_1, \dots, x_i, \dots, x_{n_0}]^T$

Matrix-multiply version: Given

W $n_1 \times n_0$

\mathbf{b} an $n_1 \times 1$ bias vector,

\mathbf{x} $n_0 \times 1$

g an activation function that is applied elementwise to a vector.

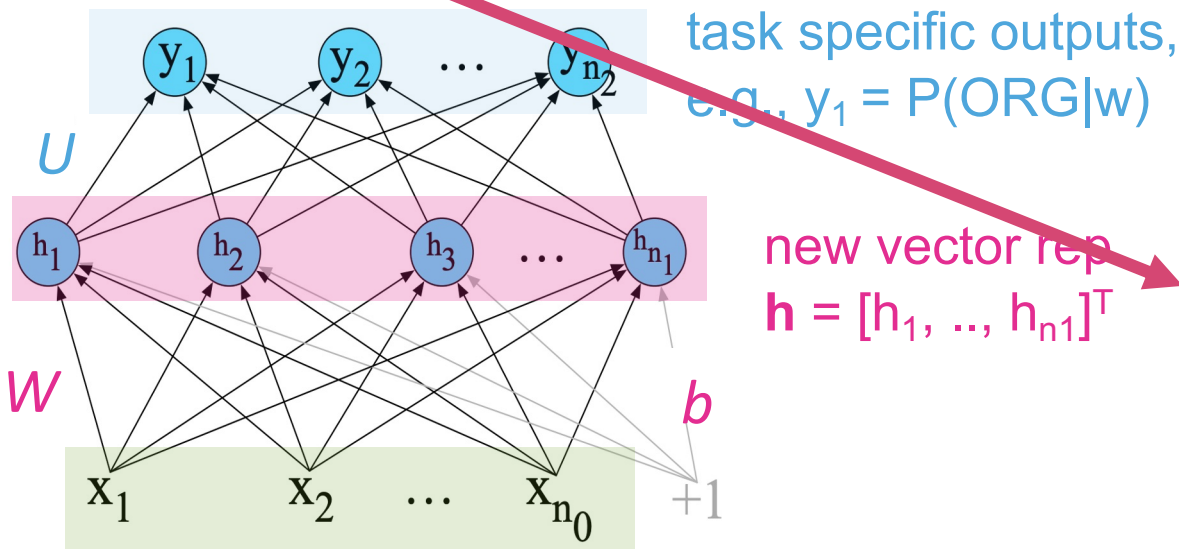
$\mathbf{w}^{[i]}$:= row i of W is the weights applied to the x_j s to help produce h_i .

$h_i = g(\mathbf{w}^{[i]} \mathbf{x} + b_i)$ or

$\mathbf{h} = g(W \mathbf{x} + \mathbf{b})$

$y_j = \text{softmax}(\mathbf{u}^{[j]} \mathbf{h})$ or $\mathbf{y} = \text{softmax}(U \mathbf{h})$

Here, \mathbf{x} is a column vector; the weights for computing h_i are row vectors.



task specific outputs, e.g., $y_1 = P(\text{ORG}|\mathbf{w})$

new vector rep
 $\mathbf{h} = [h_1, \dots, h_{n_1}]^T$

a word2vec input

$$\mathbf{x} = [x_1, \dots, x_i, \dots, x_{n_0}]^T$$

Matrix-multiply version: Given

$$W \ n_1 \times n_0$$

\mathbf{b} an $n_1 \times 1$ bias vector,

$$\mathbf{x} \ n_0 \times 1$$

g an activation function that is applied elementwise to a vector.

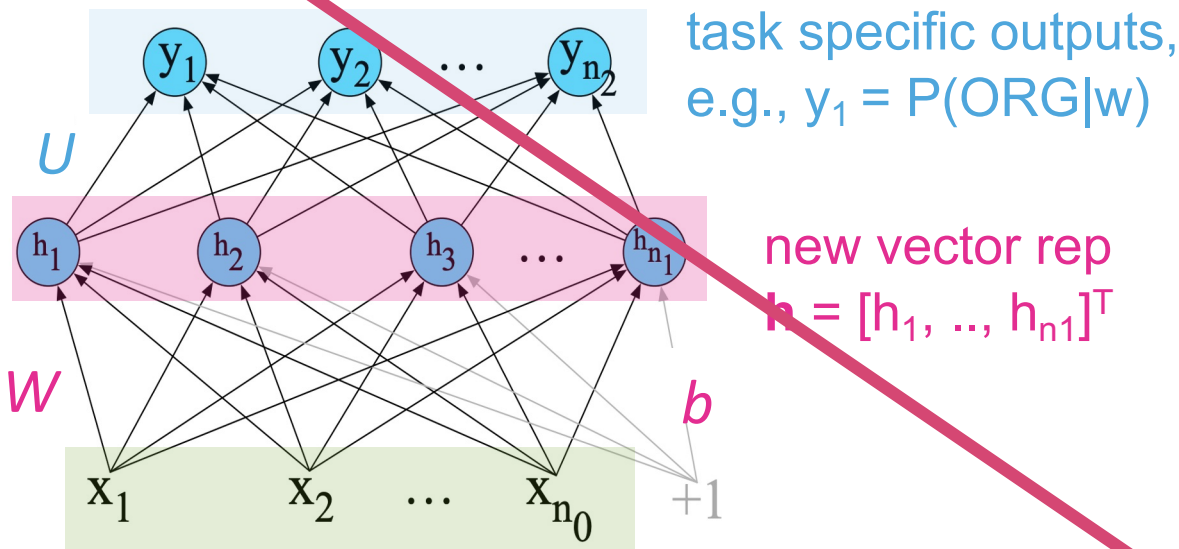
$\mathbf{w}^{[i]}$:= row i of W is the weights applied to the x_j s to help produce h_i .

$$h_i = g(\mathbf{w}^{[i]} \mathbf{x} + b_i) \text{ or}$$

$$\mathbf{h} = g(W \mathbf{x} + \mathbf{b})$$

$$y_j = \text{softmax}(\mathbf{u}^{[j]} \mathbf{h}) \text{ or } \mathbf{y} = \text{softmax}(U \mathbf{h})$$

Here, \mathbf{h} is a column vector; the weights for computing y_j are row vectors.



task specific outputs, e.g., $y_1 = P(\text{ORG}|\mathbf{w})$

new vector rep
 $\mathbf{h} = [h_1, \dots, h_{n_1}]^T$

word2vec vector as input
 $\mathbf{x} = [x_1, \dots, x_i, \dots, x_{n_0}]^T$

Matrix-multiply version: Given

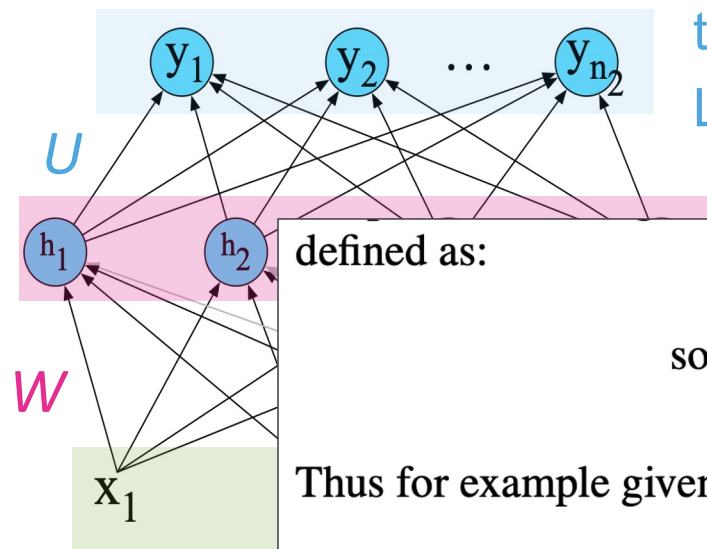
W $n_1 \times n_0$
 \mathbf{b} an $n_1 \times 1$ bias vector,
 \mathbf{x} $n_0 \times 1$
 g an activation function that is applied elementwise to a vector.

$\mathbf{w}^{[i]}$:= row i of W is the weights applied to the x_j s to help produce h_i .

$h_i = g(\mathbf{w}^{[i]} \mathbf{x} + b_i)$ or
 $\mathbf{h} = g(W \mathbf{x} + \mathbf{b})$

$y_j = \text{softmax}(\mathbf{u}^{[j]} \mathbf{h})$ or $\mathbf{y} = \text{softmax}(U \mathbf{h})$

From \mathbf{h} to $\mathbf{y} = \text{softmax}(\mathbf{U}\mathbf{h})$: how do we make the y_i 's into a probability distribution?



task specific outputs, e.g., $y_1 = P(\text{ORG}|w)$.
Let \mathbf{z} be $\mathbf{U}\mathbf{h}$. Then, $y_i = \text{softmax}_z(z_i)$.

defined as:

$$\text{softmax}(\mathbf{z}_i) = \frac{\exp(\mathbf{z}_i)}{\sum_{j=1}^d \exp(\mathbf{z}_j)} \quad 1 \leq i \leq d \quad (7.9)$$

Thus for example given a vector

$$\mathbf{z} = [0.6, 1.1, -1.5, 1.2, 3.2, -1.1], \quad (7.10)$$

the softmax function will normalize it to a probability distribution (shown rounded):

$$\text{softmax}(\mathbf{z}) = [0.055, 0.090, 0.0067, 0.10, 0.74, 0.010] \quad (7.11)$$

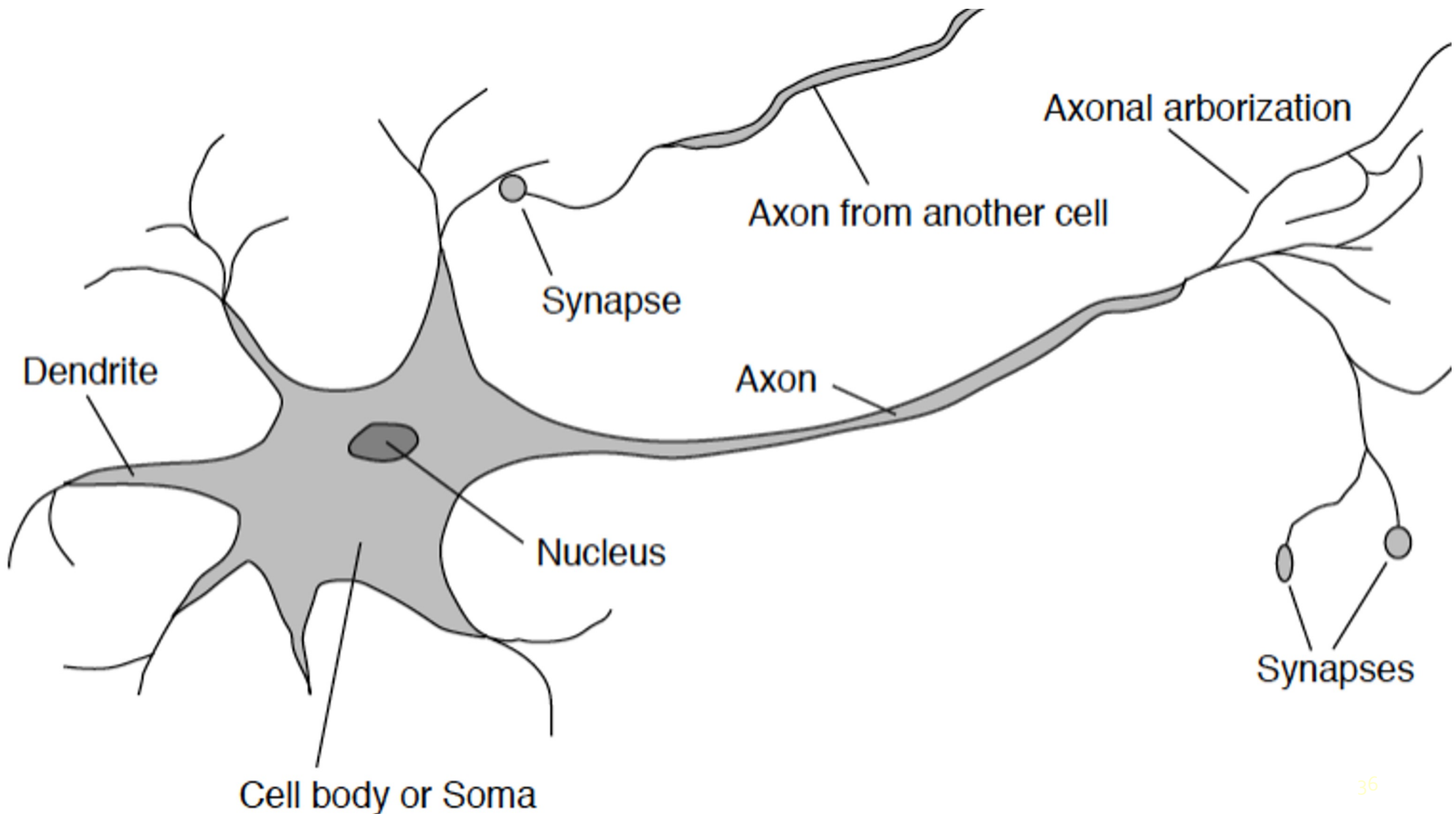
for word w ,
 $\mathbf{x} = [x_1, \dots, x_n]$

Historical note

- Rich history, starting in the early 1940's
 - McCulloch and Pitts, 1943
- Two views of neural nets
 - Models of the human brain

 Representations of complex functions

Much progress on both fronts!!!



So where do the weights come from?

Next class...