# NEURAL LANGUAGE MODELS

# Backgrounds

Computation ↔ Function
Computational Intelligence ↔ Learning the function

# Function
# $f: M \rightarrow N$

Transferring elements of source space (M) to target space (N)

$$f: R \rightarrow R : f(x) = y = x \times 2 + 1$$

$f$ is known

$f$ has two parameters (a=2, b=1)

# Function
## $f: M \rightarrow N$

Transferring elements of source space (M) to target space (N)

$f: R \rightarrow R : \{(x, f(x))| (1,1), (2,4), (3,9), ...\}$

$f(x) = ?$

$f$ is unknown!

# Function
## $f: M \rightarrow N$

Transferring elements of source space (M) to target space (N)

$f: R \rightarrow R : \{(x, f(x))| (1,1), (2,4), (3,9), ...\}$

$f(x) = ?$

$f$ is unknown! Might not exist!

# Function
## $f: M \rightarrow N$

Transferring elements of source space (M) to target space (N)

$f: R \rightarrow R : \{(x, f(x)) | (1,1), (2,4), (3,9), ...\}$

$f(x) = ?$

$f$ is unknown! If exists, there may be more that one!

# Function
## $f: M \rightarrow N$

Transferring elements of source space (M) to target space (N)
$f: R \rightarrow R : \{(x, f(x))| (1,1), (2,4), (3,9), ...\}$
$f(x) = ?$

$f$ is unknown! Even more than one exists, we may not find one!

# Function
## $f: M \longrightarrow N$

Transferring elements of source space (M) to target space (N)
$f: R \longrightarrow R : \{(x, f(x))| (1,1), (2,4), (3,9), ...\}$
$f(x) = ?$
$f$ is unknown! We assume at least one exists, we try to find it!

# Transformation
## $f: M \rightarrow N$

Transferring points (vectors) of source space (M) to target space (N)

$[T][X] = [Y]; X \in M, Y \in N$

T is a matrix that includes the parameters of $f$

# Transformation (Linear Algebra)
## $f: M \rightarrow N$

Transferring points (vectors) of source space (M) to target space (N)

$$[2\ 1]\begin{bmatrix} x \\ 1 \end{bmatrix} = [2x + 1]$$

$$f(x) = 2x + 1$$

# Transformation (Linear Algebra)
## $f: M \rightarrow N$

Transferring points (vectors) of source space (M) to target space (N)

$(AB)^T = B^T A^T$ →

$$[x \ 1]\begin{bmatrix} 2 \\ 1 \end{bmatrix} = [2x + 1]$$

$f(x) = 2x + 1$

# Transformation (Linear Algebra)
## $f: M \longrightarrow N$

Transferring points (vectors) of source space (M) to target space (N)

$$[x \ 1]\begin{bmatrix} \mathbf{w_1} \\ \mathbf{w_2} \end{bmatrix} = \{(x, f(x))|\ (1,1),\ (2,4),\ (3,9),\ ...\}$$

$$f(x) = ?$$

# Neural Network
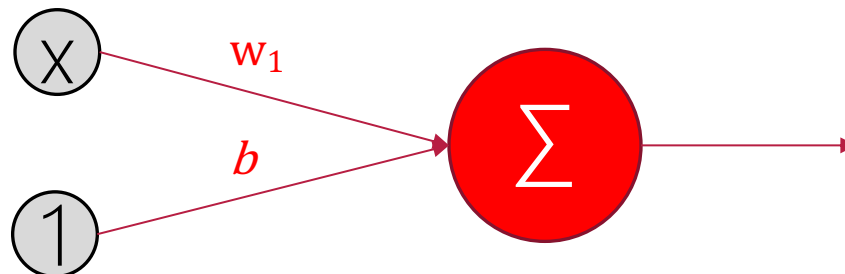## *f: M → N*
## *Perceptron*

Transferring points (vectors) of source space (M) to target space (N)

$[x \ 1] \begin{bmatrix} \mathbf{w_1} \\ \boldsymbol{b} \end{bmatrix} = y$
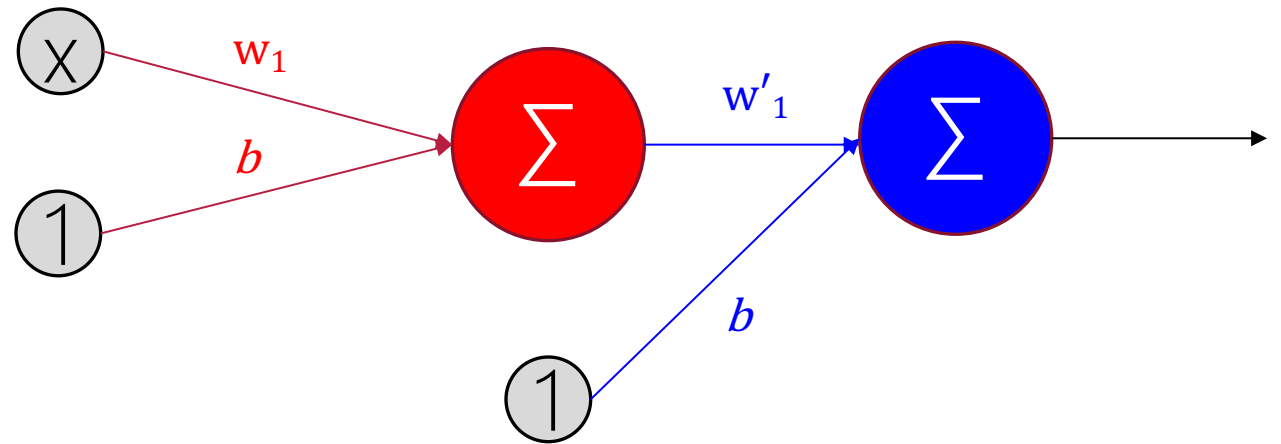
$f(x) = y$

# Neural Network

$$g: (f: M \rightarrow N) \rightarrow O$$

Transferring points (vectors) of source space (M) to target space (O)

$$[[x\ 1]\begin{bmatrix} \mathbf{w}_1 \\ \boldsymbol{b} \end{bmatrix}\ 1]\begin{bmatrix} \mathbf{w'}_1 \\ \boldsymbol{b} \end{bmatrix} = z$$

$$(g \circ f)(x) = g(f(x)) = z$$

# Neural Network

$$g: (f: M \rightarrow N) \rightarrow O$$

Transferring points (vectors) of source space (M) to target space (O)

$$[y \ 1]\begin{bmatrix} \mathbf{w'}_1 \\ b \end{bmatrix} = z$$
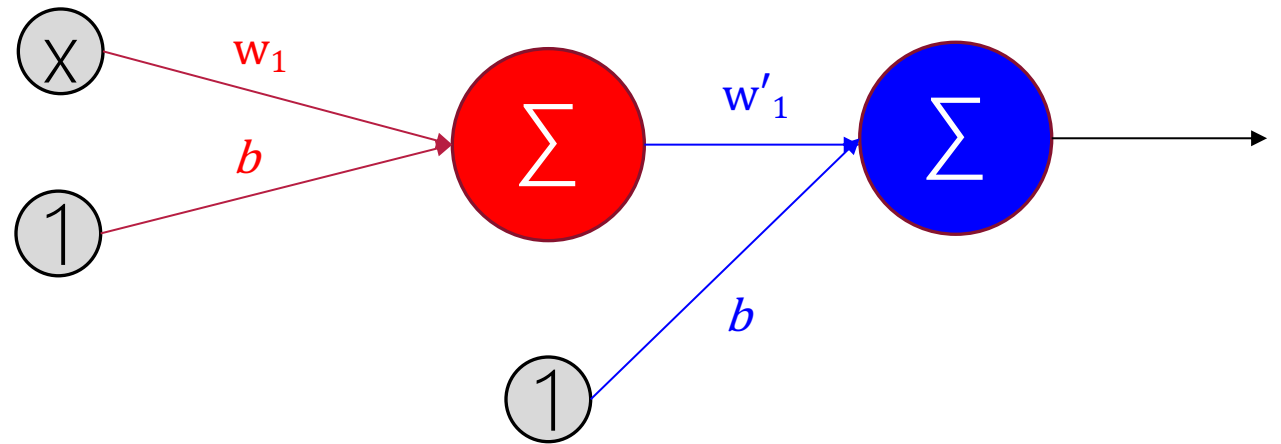
$$g(f(x)) = g(y) = z$$

# Neural Network

$$g: (f: M \rightarrow N) \rightarrow O$$

Transferring points (vectors) of source space (M) to target space (O)

$$[y \ 1]\begin{bmatrix} 0 \\ \sigma(\mathbf{y}) \end{bmatrix} = y$$

$$g(f(x)) = g(y) = z$$



$$\left[[x \ 1]\begin{bmatrix} w_1 \\ b \end{bmatrix} \quad 1\right]\begin{bmatrix} 0 \\ \sigma(x \times w_1 + b) \end{bmatrix} = z$$
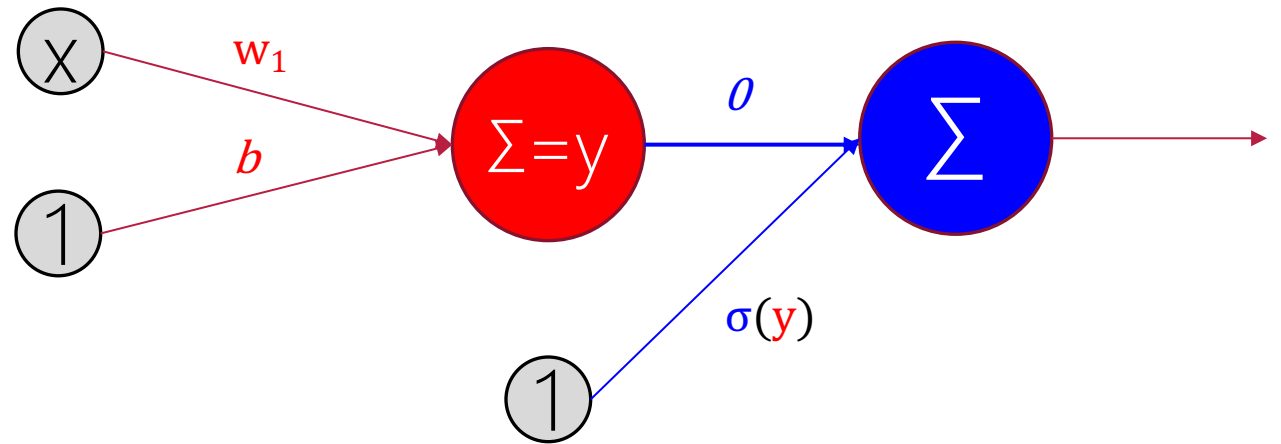
# Neural Network
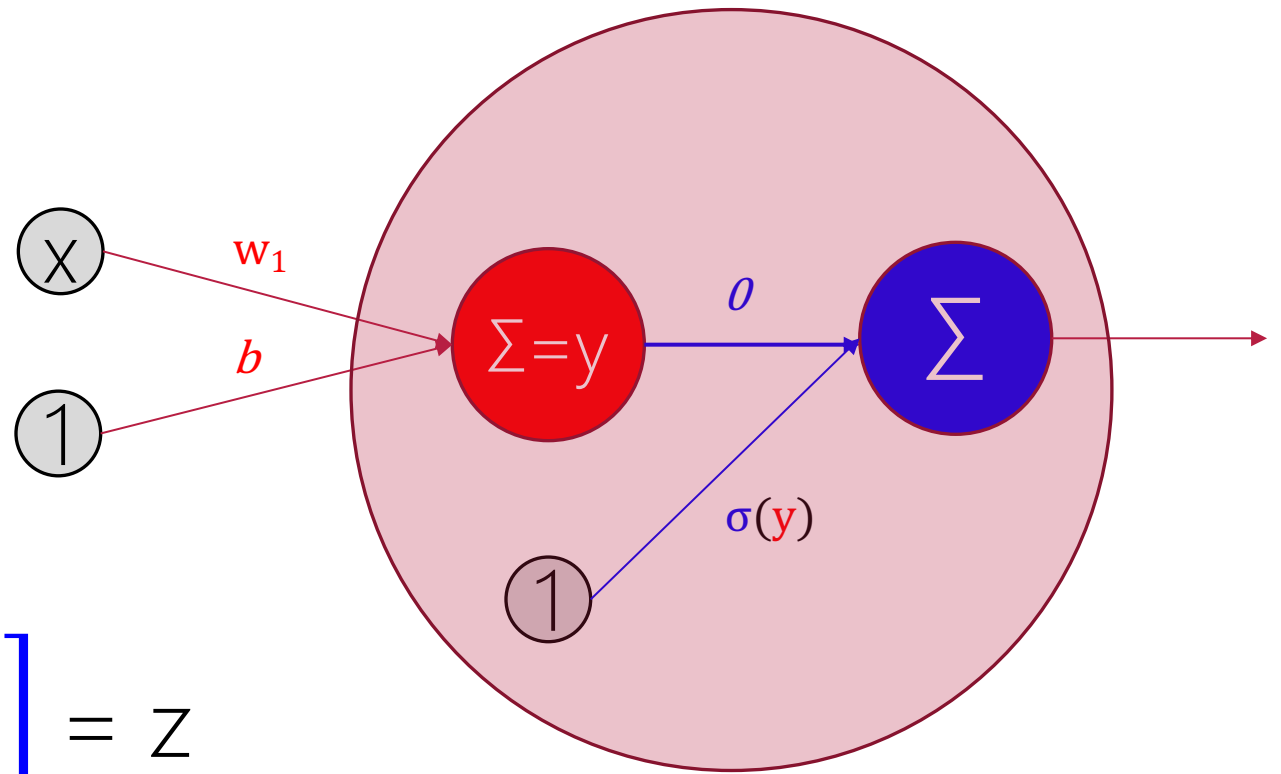
$$g: (f: M \rightarrow N) \rightarrow O$$

Transferring points (vectors) of source space (M) to target space (O)

$$[y \ 1]\begin{bmatrix} 0 \\ \sigma(\mathbf{y}) \end{bmatrix} = y$$

$$g(f(x)) = g(y) = z$$



$$\left[[x \ 1]\begin{bmatrix} w_1 \\ b \end{bmatrix} \ 1\right]\begin{bmatrix} 0 \\ \sigma(x \times w_1 + b) \end{bmatrix} = z$$
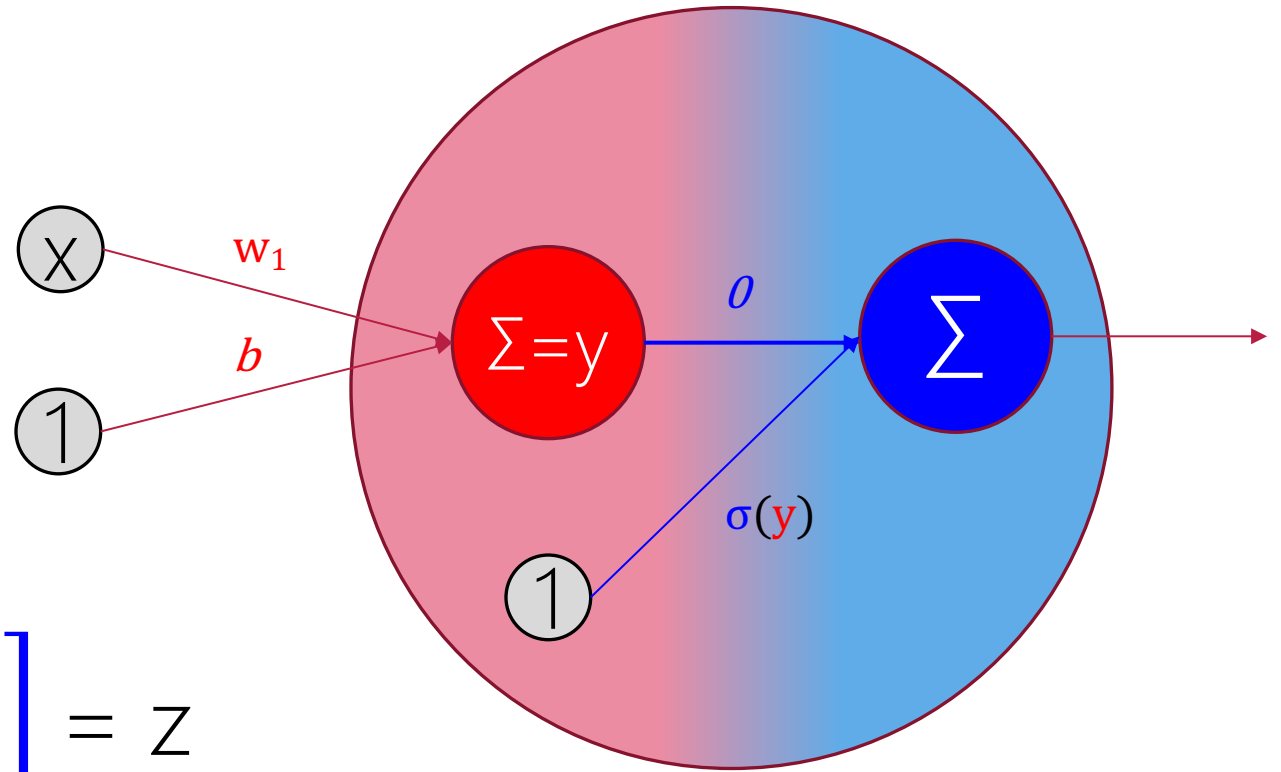
# Neural Network

$$h: M \longrightarrow O$$

Transferring points (vectors) of source space (M) to target space (O)

$$h(x) = \sigma(x \times w_1 + b)$$



$$[[x\ 1]\begin{bmatrix} w_1 \\ b \end{bmatrix}\ \ 1]\begin{bmatrix} 0 \\ \sigma(x \times w_1 + b) \end{bmatrix} = z$$

# Neural Network

$$h: M \rightarrow O$$

Transferring points (vectors) of source space (M) to target space (O)

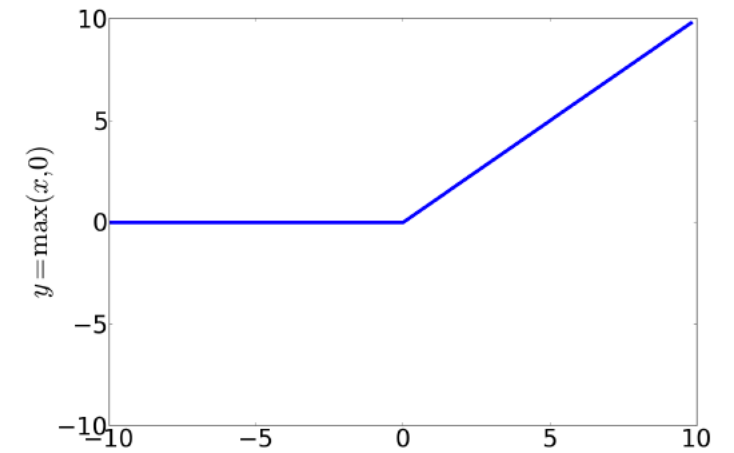h(x) = $\sigma$(x × $w_1$+ $b$)

$\sigma$ is a known non-linear function: Sigmoid, ReLU, Tanh

$$[[x\ 1]\begin{bmatrix} w_1 \\ b \end{bmatrix}\ \ 1]\begin{bmatrix} 0 \\ \sigma(x \times w_1 + b) \end{bmatrix} = z$$

# Neural Network: Activation Function σ

# Neural Network

$$h: M=\{0,1\}^2 \longrightarrow N=\{0,1\}$$

M = {(0,0), (0,1), (1,0), (1,1)}
N = {0, 1}
h = {((0,0), 0), ((0,1), 0), ((1,0), 0), ((1,1), 1)}

$$[[x_1 \ \ x_2 \ \ 1]\begin{bmatrix} w_1 \\ w_2 \\ b \end{bmatrix} 1] \left[ \sigma(x_1 \times w_1 + x_2 \times w_2 + b) \right] = z$$

# Neural Network

## *AND: M={0,1}² → N={0,1}*

M = {(0,0), (0,1), (1,0), (1,1)}

N = {0, 1}

h = {((0,0), 0), ((0,1), 0), ((1,0), 0), ((1,1), 1)}



$$[[x_1 \quad x_2 \quad 1]\begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix} 1] [max(0, x_1 \times 1 + x_2 \times 1 + -1)] = z$$

# Neural Network

*OR: M={0,1}² → N={0,1}*

M = {(0,0), (0,1), (1,0), (1,1)}

N = {0, 1}

h = {((0,0), 0), ((0,1), 1), ((1,0), 1), ((1,1), 1)}

$$[[x_1 \ x_2 \ 1] \begin{bmatrix} \textbf{w}_1 \\ \textbf{w}_2 \\ \textbf{\textit{b}} \end{bmatrix} \ 1] \begin{bmatrix} \textit{0} \\ \sigma(x_1 \times \textbf{w}_1 + x_2 \times \textbf{w}_2 + \textbf{\textit{b}}) \end{bmatrix} = z$$

# Neural Network

## *XOR: M={0,1}² → N={0,1}*

M = {(0,0), (0,1), (1,0), (1,1)}

N = {0, 1}

h = {((0,0), 0), ((0,1), 1), ((1,0), 1), ((1,1), 0)}

$$[[x_1 \ x_2 \ 1] \begin{bmatrix} \mathbf{w_1} \\ \mathbf{w_2} \\ \boldsymbol{b} \end{bmatrix} \ 1] \begin{bmatrix} 0 \\ \sigma(x_1 \times \mathbf{w_1} + x_2 \times \mathbf{w_2} + \boldsymbol{b}) \end{bmatrix} = z$$

# Neural Network

*XOR: M={0,1}$^2$ → N={0,1,2}$^2$ → O={0,1}*

M = {(0,0), (0,1), (1,0), (1,1)} → N = {(0,0), (0,1), ..., (2,2)} → O = {0, 1}

h = {((0,0), 0), ((0,1), 1), ((1,0), 1), ((1,1), 0)}

$$[[x_1 \quad x_2 \quad 1] \begin{bmatrix} w_1 \\ w_2 \\ b \end{bmatrix} \quad 1] \begin{bmatrix} 0 \\ \sigma(x_1 \times w_1 + x_2 \times w_2 + b) \end{bmatrix} = y_1$$
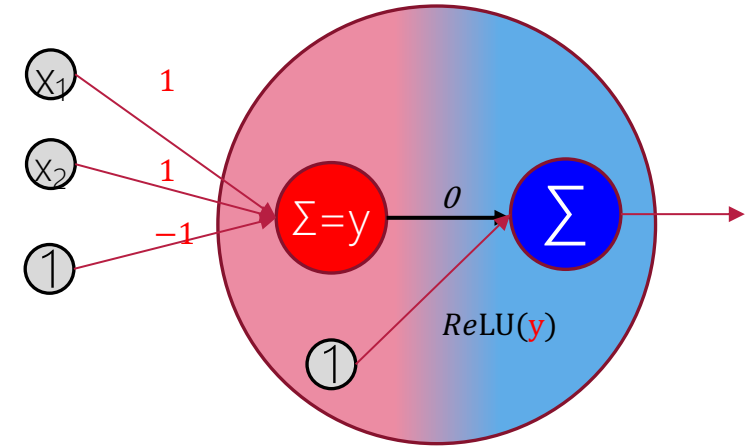
$$[[x_1 \quad x_2 \quad 1] \begin{bmatrix} w'_1 \\ w'_2 \\ b' \end{bmatrix} \quad 1] \begin{bmatrix} 0 \\ \sigma(x_1 \times w'_1 + x_2 \times w'_2 + b') \end{bmatrix} = y_2$$

$$[[y_1 \quad y_2 \quad 1] \begin{bmatrix} w''_1 \\ w''_2 \\ b'' \end{bmatrix} \quad 1] \begin{bmatrix} 0 \\ \sigma(x_1 \times w''1 + x_2 \times w''2 + b'') \end{bmatrix} = z$$

# Neural Network

*XOR: M={0,1}² → N={0,1,2}² → O={0,1}*

M = {(0,0), (0,1), (1,0), (1,1)} → N = {(0,0), (0,1), ..., (2,2)} → O = {0, 1}

h = {((0,0), 0), ((0,1), 1), ((1,0), 1), ((1,1), 0)}

$$[[x_1 \ x_2 \ 1] \begin{bmatrix} w_1 & w'1 \\ w_2 & w'2 \\ b & b' \end{bmatrix} \ 1] =$$

$$[x_1 \times w_1 + x_2 \times w_2 + b \quad x_1 \times w'1 + x_2 \times w'2 + b' \quad 1] \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ \sigma(x_1 \times w_1 + x_2 \times w_2 + b) & \sigma(x_1 \times w'1 + x_2 \times w'2 + b') \end{bmatrix} =$$

$$[\underbrace{\sigma(x_1 \times w_1 + x_2 \times w_2 + b)}_{\alpha} \ \underbrace{\sigma(x_1 \times w'1 + x_2 \times w'2 + b')}_{\beta} \ 1] \begin{bmatrix} w''_1 \\ w''_2 \\ b'' \end{bmatrix} \ 1] \begin{bmatrix} 0 \\ 0 \\ \sigma(\alpha \times w''1 + \beta \times w''2 + b'') \end{bmatrix} = z$$
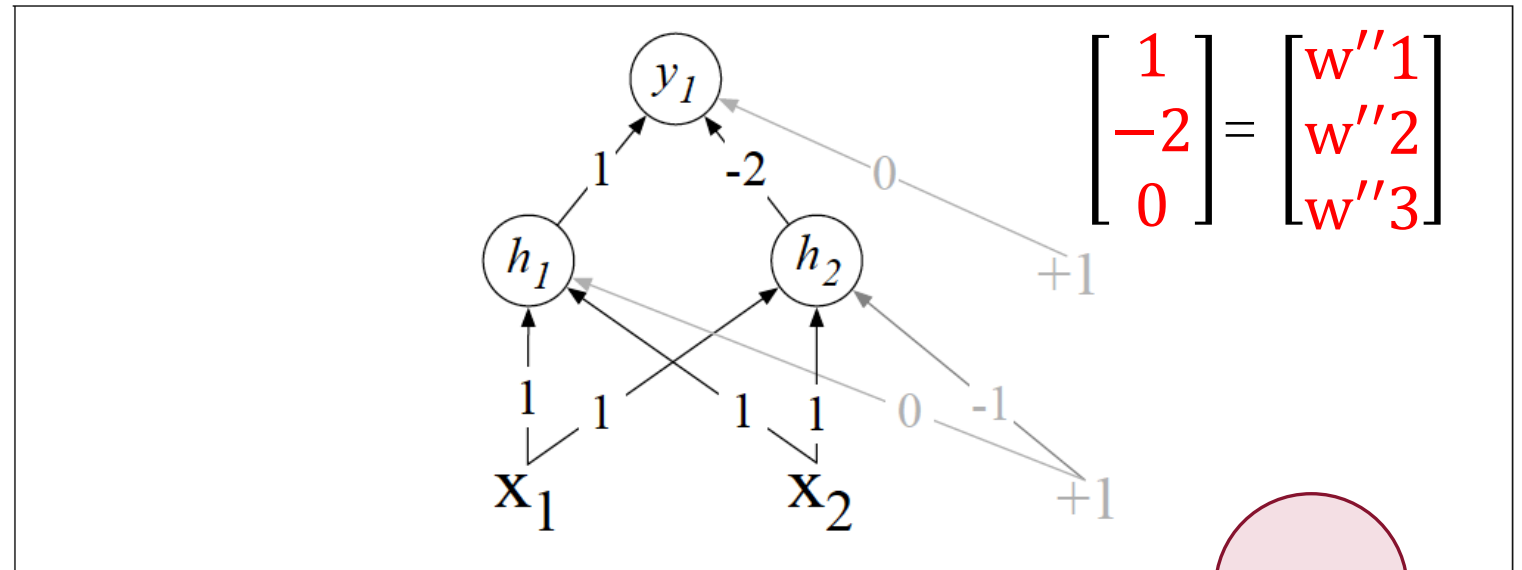
# Neural Network

## XOR: $M=\{0,1\}^2 \rightarrow N=\{0,1,2\}^2 \rightarrow O=\{0,1\}$

$M = \{(0,0), (0,1), (1,0), (1,1)\} \rightarrow N = \{(0,0), (0,1), ..., (2,2)\} \rightarrow O = \{0, 1\}$

$h = \{((0,0), 0), ((0,1), 1), ((1,0), 1), ((1,1), 0)\}$

$$[x_1 \ x_2 \ 1]\begin{bmatrix} w_1 & w'1 \\ w_2 & w'2 \\ w_3 & w'3 \end{bmatrix} =$$

$$[x_1 \ x_2 \ 1]\begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 0 & -1 \end{bmatrix} =$$

$[h_1=\text{ReLU}(x_1+x_2) \quad h_2=\text{ReLU}(x_1+x_2-1)]$

$$\begin{bmatrix} 1 \\ -2 \\ 0 \end{bmatrix} = \begin{bmatrix} w''1 \\ w''2 \\ w''3 \end{bmatrix}$$



**Figure 7.6**  XOR solution after Goodfellow et al. (2016). There are three ReLU units, in two layers; we've called them $h_1$, $h_2$ ($h$ for "hidden layer") and $y_1$. As before, the numbers on the arrows represent the weights $w$ for each unit, and we represent the bias $b$ as a weight on a unit clamped to +1, with the bias weights/units in gray.

# Neural Network

## XOR: $M=\{0,1\}^2 \rightarrow N=\{0,1,2\}^2 \rightarrow O=\{0,1\}$

M = {(0,0), (0,1), (1,0), (1,1)} → N = {(0,0), (0,1), ..., (2,2)} → O = {0, 1}

$(0,0) \rightarrow (0,0) \rightarrow (0)$
$(0,1) \rightarrow (1,0) \rightarrow (1)$
$(1,0) \rightarrow (1,0) \rightarrow (1)$
$(1,1) \rightarrow (2,1) \rightarrow (0)$



a) The original $x$ space          b) The new $h$ space

**Figure 7.7**   The hidden layer forming a new representation of the input. Here is the representation of the hidden layer, $h$, compared to the original input representation $x$. Notice that the input point [0 1] has been collapsed with the input point [1 0], making it possible to linearly separate the positive and negative cases of XOR. After Goodfellow et al. (2016).

# Universal Approximation Theorem
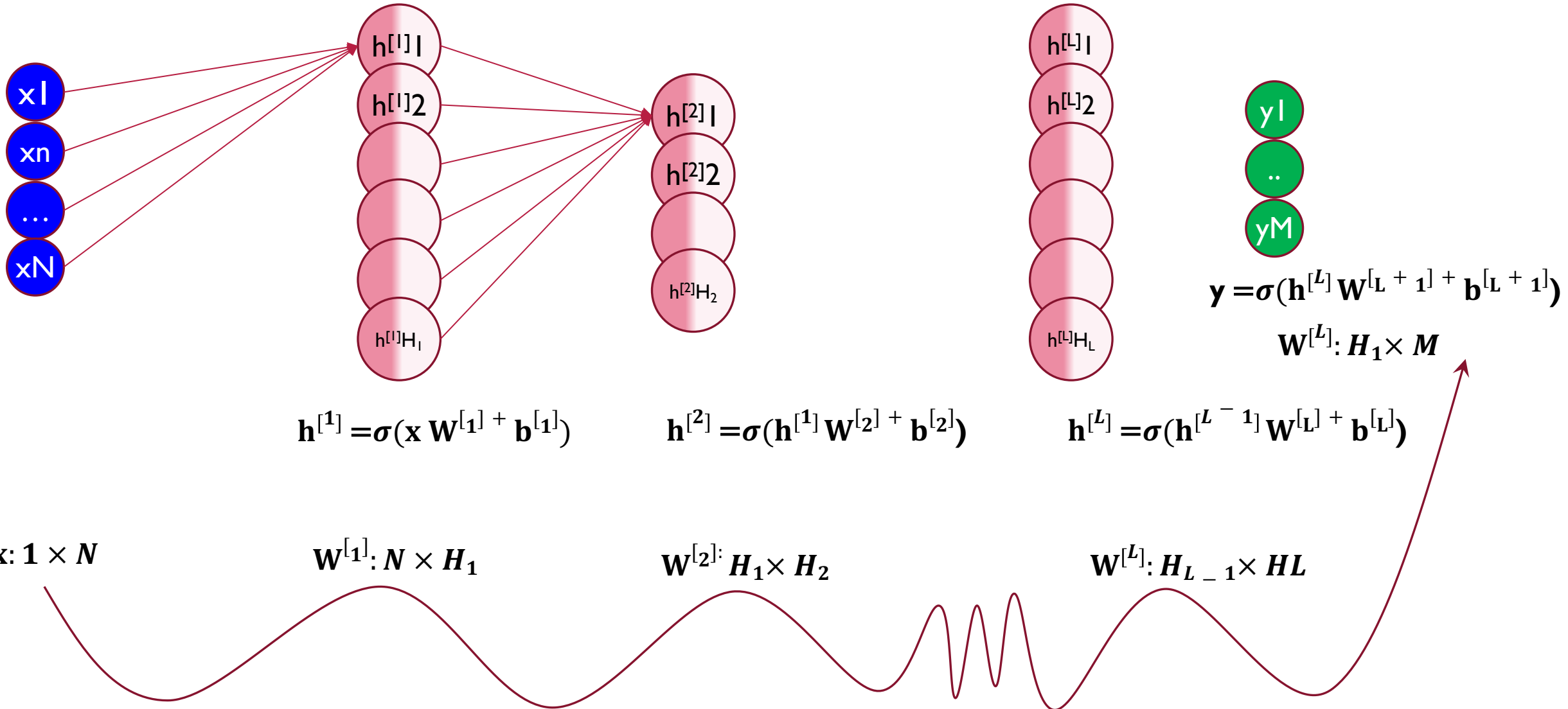
Be carful, it shows the existence (power) of neural nets, but it does not show which architecture is the function!

How many transformation (#layers)?
To what intermediate space (#nodes)?

# Neural Network: Feed-forward



$$\mathbf{h}^{[1]} = \sigma(\mathbf{x}\,\mathbf{W}^{[1]} + \mathbf{b}^{[1]})$$

$$\mathbf{h}^{[2]} = \sigma(\mathbf{h}^{[1]}\,\mathbf{W}^{[2]} + \mathbf{b}^{[2]})$$

$$\mathbf{h}^{[L]} = \sigma(\mathbf{h}^{[L-1]}\,\mathbf{W}^{[L]} + \mathbf{b}^{[L]})$$

$$\mathbf{y} = \sigma(\mathbf{h}^{[L]}\,\mathbf{W}^{[L+1]} + \mathbf{b}^{[L+1]})$$

$$\mathbf{W}^{[L]}: H_1 \times M$$

$$\mathbf{x}: 1 \times N \qquad \mathbf{W}^{[1]}: N \times H_1 \qquad \mathbf{W}^{[2]}: H_1 \times H_2 \qquad \mathbf{W}^{[L]}: H_{L-1} \times HL$$

# Neural Network: Feed-forward Train

The only known is the data (input: X, output: Y)
Given the input point X, the net should land it to Y.

xl

xn

...

xN

yl

..

yM

If the net should land it to $f(X) = Y' \rightarrow$ error!
Correct the net to make Y' close to Y.

How?

# Neural Network: Feed-forward Train

x1
xn
...
xN

y1
..
yM

How?
By minimizing the overall error.
What is error?
e = distance (Y, f(X)); for all (X, Y)

Optimization Function ➔ Loss Function
Minimization

# Neural Network: Gradient

How?
By minimizing the overall error.
What is error?

xl

xn

...
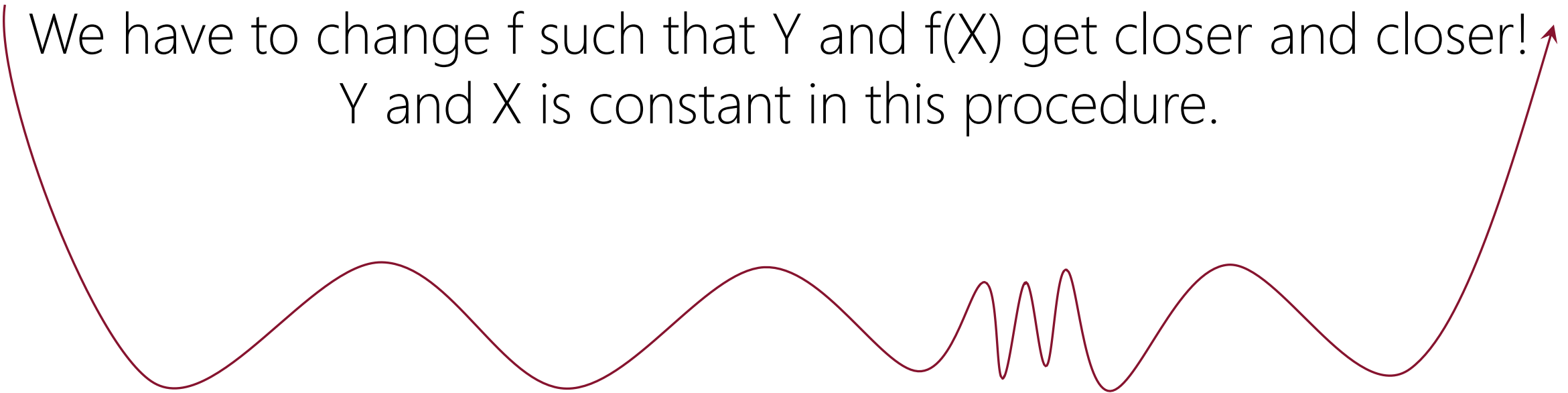
xN

yl

..

yM

e = distance (Y, f(X)); for all (X, Y)

We have to change f such that Y and f(X) get closer and closer!
Y and X is constant in this procedure.

# Neural Network: Gradient

How?
By minimizing the overall error.
What is error?

xl

xn

...

xN

yl

..

yM

e = distance (Y, f(X)); for all (X, Y)

We have to change parameters such that Y and f(X) get closer and closer!
Y and X is constant in this procedure.
$e(\Theta) = distance(Y, f(X; \Theta))$

# Neural Network: Gradient

$f(x) = 3x+1$ ➔ $f(x; \Theta) = f(x; [3,1]) = 3x+1$

$f(x; [w_1,w_2]) = w_1 x + w_2$

$(x, y) = (2, 1)$

xl

xn

...

xN

yl

..

yM

$f(2; [w_1,w_2]) = w_1 2 + w_2$

$e = e(w_1,w_2) = e(\Theta) = |f(2; [w_1,w_2]) - 1| = w_1 2 + w_2 - 1$

How to change $w_1$ and $w_2$ to reduce the error?

# Neural Network: Gradient Descent

f(x) = 3x+1 ➜ f(x; $\Theta$) = f(x; [3,1]) = 3x+1

$$f(x; [w_1,w_2]) = w_1 x + w_2$$

(x, y) = (2, 1)

xl

xn

...

xN

yl

..

yM

$$f(2; [w_1,w_2]) = w_1 2 + w_2$$

$$e = e(w_1,w_2) = e(\Theta) = |\, f(2; [w_1,w_2]) - 1\,| = w_1 2 + w_2 - 1$$

$$w_1 = w_1 - \eta \frac{\partial e}{\partial w_1} = w_1 - 0.1 \times 2 = w_1 - 0.2$$

$$w_2 = w_2 - \eta \frac{\partial e}{\partial w_2} = w_2 - 0.1 \times 1 = w_1 - 0.1$$

# Neural Network: Gradient Descent

$$f(x; \Theta) = f(x; [2.8, 0.9]) = 2.8x + 0.9 =$$
$$f(x; \Theta) = f(x; [2.6, 0.8]) = 2.6x + 0.8 =$$
$$f(x; \Theta) = f(x; [2.4, 0.7]) = 2.4x + 0.7 =$$
$$\ldots$$

x1

xn

...

xN

y1

..

yM

# Neural Network: Gradient Descent

xl

xn

...

xN

yl

..

yM

Convex Error Function
Many Parameters (all matrices' elements)

Backpropagation!
Computation Graph
Forward Pass
Backward Differentiation.

# Neural Language Model

# Neural Language Model

## Bigram Language Model

$$P(w_i | w_{i-1}) = \frac{\#(wi, wi_{-1)}}{\#(wi)}$$

$$f: (\{0,1\}^{|V|})^2 \longrightarrow R^{[0,1]}; \; f(w_i, w_{i-1}) = P(w_i | w_{i-1})$$

Input: one-hot vector                    Output:

$w_{i-1} = [0\ 0\ 0\ ...\ 1\ 0...\ 0\ 0\ ]$          [y]

$w_i\quad = [0\ 0\ 0\ ...\ 0\ 1\ ...\ 0\ 0\ ]$

# Neural Language Model

## Bigram Language Model

$$P(w_i | w_{i-1}) = \frac{\#(w_i, w_{i-1})}{\#(w_i)}$$

$$f: (R^{|V|})^2 \longrightarrow R^{[0,1]}; \; f(w_i, w_{i-1}) = P(w_i | w_{i-1})$$

Input: sparse semantic vector

$w_{i-1} = [1\ 2\ 3\ ...\ 1\ 2...\ 1\ 0\ ]$

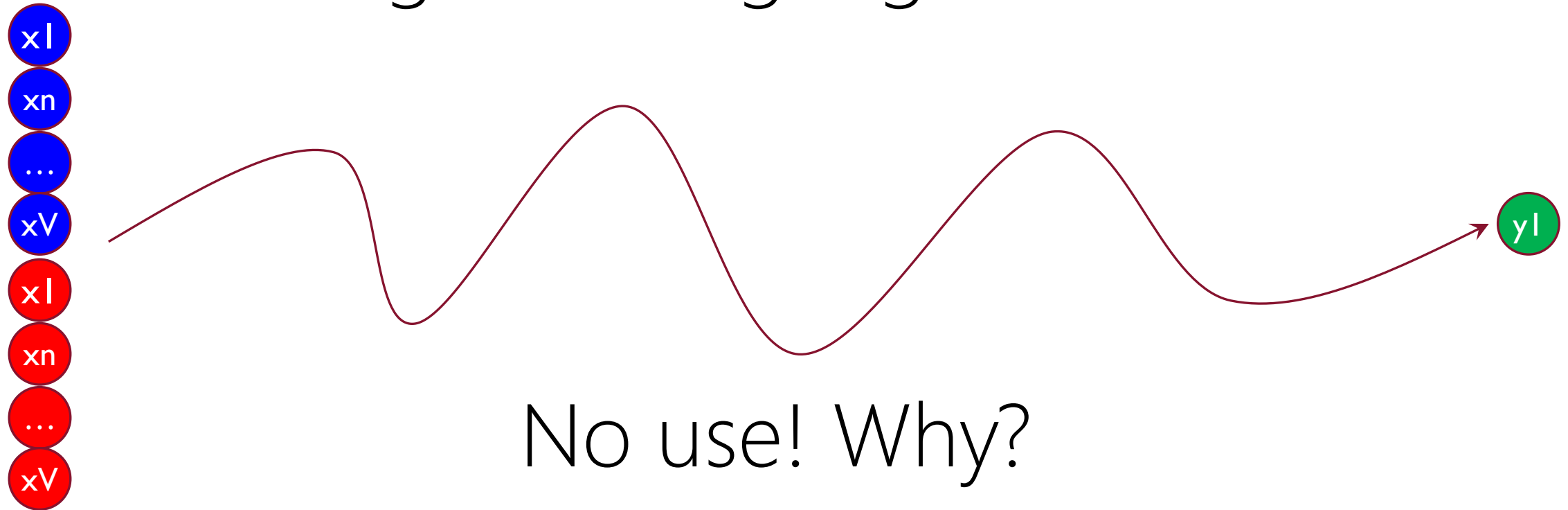$w_i\quad = [4\ 0\ 2\ ...\ 0\ 5\ ...\ 3\ 0\ ]$

Output:

$[y]$
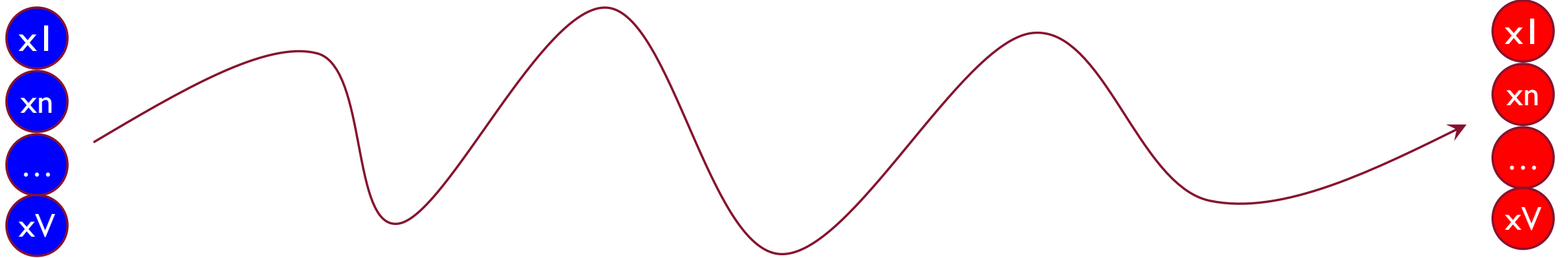
# Neural Language Model

## Bigram Language Model



No use! Why?

# Neural Language Model

## Bigram Language Model



$P(w_i| w_{i-1}) = ?$ Unknown Initially

While enumerating the text stream, when I see $[w_{i-1} \ w_i]$ then

$$f: \{0,1\}^{|V|} \rightarrow \{0,1\}^{|V|}; \ f(w_{i-1}) \rightarrow w_i$$

$$[0 \ 0 \ 0 \ ... \ 1 \ 0 \ ... \ 0 \ 0 \ 0] \rightarrow [0 \ 0 \ 0 \ ... \ 0 \ 1 \ ... \ 0 \ 0 \ 0]$$

# Neural Language Model

## Bigram Language Model



P($w_i$| $w_{i-1}$) = ? Unknown Initially
While enumerating the text stream, when I see [$w_{i-1}$ $w_i$] then
f: $\{0,1\}^{|V|} \rightarrow R^{|V|}$; f($w_{i-1}$) $\rightarrow$ #($w_{i-1}w_i$)
[0 0 0 … 1 0 … 0 0 0] $\rightarrow$ [0 1 0 … 0 2 … 1 0 0]

# Neural Language Model

## Bigram Language Model



$P(w_i | w_{i-1}) = ?$ Unknown Initially

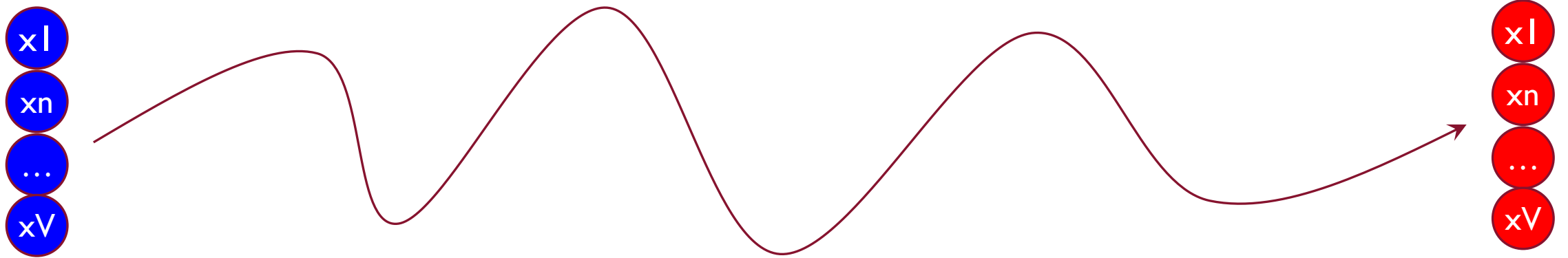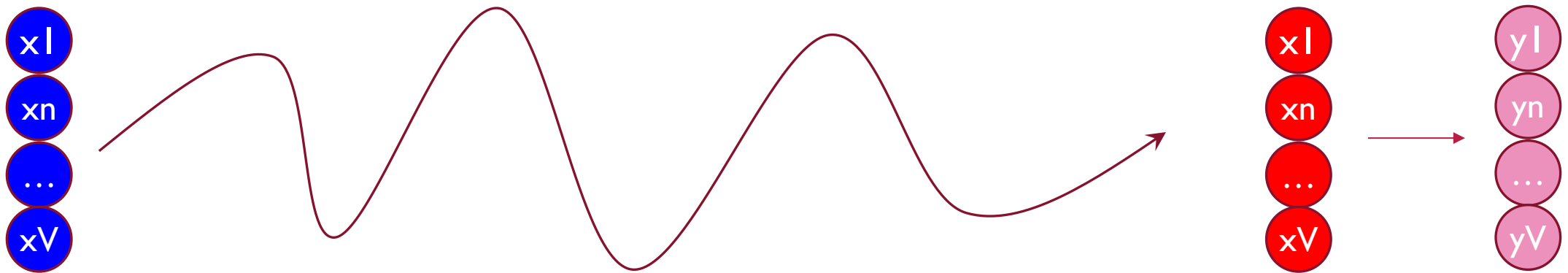While enumerating the text stream, when I see $[w_{i-1} \, w_i]$ then

$f: \{0,1\}^{|V|} \to R^{|V|}; \; f(w_{i-1}) \to \#(w_{i-1}w_i) \to$ Normalized

$[0 \; 0 \; 0 \; \dots \; 1 \; 0 \; \dots \; 0 \; 0 \; 0] \to [0 \; 1 \; 0 \; \dots \; 0 \; 2 \; \dots \; 1 \; 0 \; 0] \to$ Softmax

# Neural Language Model

## Bigram Language Model



$$y_i = Softmax(x_i) = \frac{e^{x_i}}{\sum_{j=1}^{V} e^{x_j}}$$

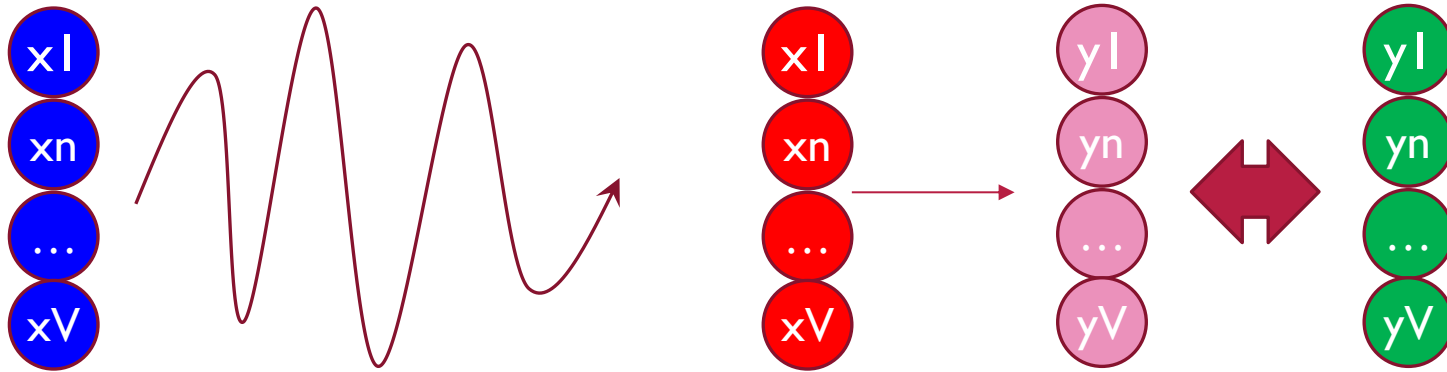$$0 \leq y_i \leq 1$$

$$\sum_{j=1}^{V} y_i = 1$$

$$f: f: \{0,1\}^{|V|} \longrightarrow [0,1]^{|V|}; f(w_{i-1}) = P(w_i | w_{i-1})$$

# Neural Language Model

## Bigram Language Model



$$\text{Error} = \text{Loss} = -\text{Log } P(w_i | w_{i-1})$$

# Neural Language Model

Bengio, Yoshua, et al. "A neural probabilistic language model." The journal of machine learning research 3 (2003): 1137-1155.



**Output layer P(w|u)** $1 \times |V|$

$|V| \times d_h \, U$

**Hidden layer** $1 \times d_h$

$d_h \times 3d \, W$

**Projection layer** $1 \times 3d$

concatenated embeddings for context words

embedding for word 35

embedding for word 9925

embedding for word 45180

$P(w_t = V_{42} | w_{t-3}, w_{t-2}, w_{t-3})$

word 42

... hole | in | the | ground | there | lived | ...

$w_{t-3}$  $w_{t-2}$  $w_{t-1}$  $w_t$

N-gram or Bag-of-Word?

# Neural Language Model

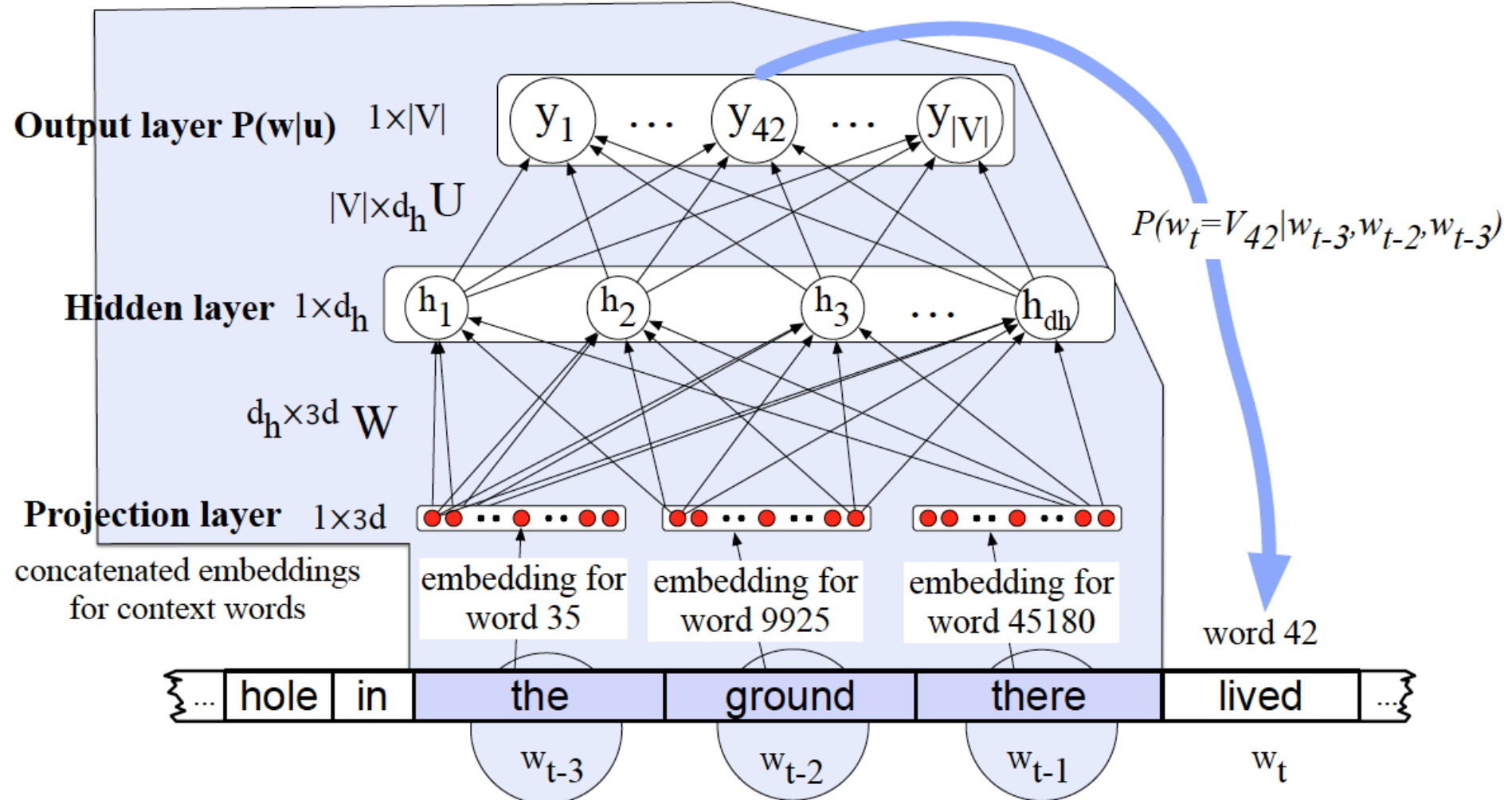Bengio, Yoshua, et al. "A neural probabilistic language model." The journal of machine learning research 3 (2003): 1137-1155.



**Output layer P(w|u)** $1 \times |V|$ — $y_1 \cdots y_{42} \cdots y_{|V|}$

$|V| \times d_h$ $U$

$P(w_t = V_{42} | w_{t-3}, w_{t-2}, w_{t-3})$

**Hidden layer** $1 \times d_h$ — $h_1 \quad h_2 \quad h_3 \cdots h_{dh}$

$d_h \times 3d$ $W$

**Projection layer** $1 \times 3d$

concatenated embeddings for context words

embedding for word 35

embedding for word 9925

embedding for word 45180

word 42

... hole in | the | ground | there | lived ...

$w_{t-3}$ $w_{t-2}$ $w_{t-1}$ $w_t$

Bottleneck → Softmax → Hierarchical Softmax

# Neural Language Model