# Introduction to neural networks

Herman Kamper

http://www.kamperh.com/

# Introduction to neural networks

## Preliminaries

Herman Kamper

http://www.kamperh.com/

# Vector and matrix derivatives recap  (Denominator layout)

$$\frac{\partial J}{\partial W}$$

Derivative of a scalar function $f : \mathbb{R}^N \to \mathbb{R}$ with respect to vector $\mathbf{x} \in \mathbb{R}^N$:

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \triangleq \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_N} \end{bmatrix}$$

Derivative of a vector function $\boldsymbol{f} : \mathbb{R}^N \to \mathbb{R}^M$, where
$\boldsymbol{f}(\mathbf{x}) = \begin{bmatrix} f_1(\mathbf{x}) & f_2(\mathbf{x}) & \cdots & f_M(\mathbf{x}) \end{bmatrix}^\top$, with respect to vector $\mathbf{x} \in \mathbb{R}^N$:

$$\frac{\partial \boldsymbol{f}(\mathbf{x})}{\partial \mathbf{x}} \triangleq \begin{bmatrix} \frac{\partial \boldsymbol{f}(\mathbf{x})}{\partial x_1} \\ \frac{\partial \boldsymbol{f}(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial \boldsymbol{f}(\mathbf{x})}{\partial x_N} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \frac{\partial f_2(\mathbf{x})}{\partial x_1} & \cdots & \frac{\partial f_M(\mathbf{x})}{\partial x_1} \\ \frac{\partial f_1(\mathbf{x})}{\partial x_2} & \frac{\partial f_2(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_M(\mathbf{x})}{\partial x_2} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial f_1(\mathbf{x})}{\partial x_N} & \frac{\partial f_2(\mathbf{x})}{\partial x_N} & \cdots & \frac{\partial f_M(\mathbf{x})}{\partial x_N} \end{bmatrix}$$
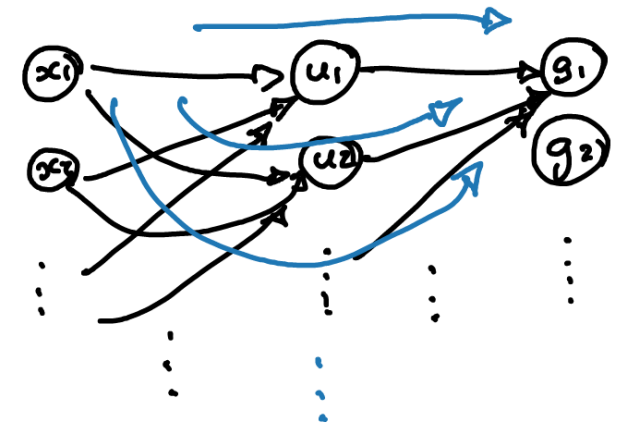
Transpose of Jacobian

Derivative of a scalar function $f : \mathbb{R}^{M \times N} \to \mathbb{R}$
with respect to matrix $\mathbf{X} \in \mathbb{R}^{M \times N}$:

$$\frac{\partial f(\mathbf{X})}{\partial \mathbf{X}} \triangleq \begin{bmatrix} \dfrac{\partial f(\mathbf{X})}{\partial X_{1,1}} & \dfrac{\partial f(\mathbf{X})}{\partial X_{1,2}} & \cdots & \dfrac{\partial f(\mathbf{X})}{\partial X_{1,N}} \\[2ex] \dfrac{\partial f(\mathbf{X})}{\partial X_{2,1}} & \dfrac{\partial f(\mathbf{X})}{\partial X_{2,2}} & \cdots & \dfrac{\partial f(\mathbf{X})}{\partial X_{2,N}} \\[2ex] \vdots & \vdots & \vdots & \vdots \\[2ex] \dfrac{\partial f(\mathbf{X})}{\partial X_{M,1}} & \dfrac{\partial f(\mathbf{X})}{\partial X_{M,2}} & \cdots & \dfrac{\partial f(\mathbf{X})}{\partial X_{M,N}} \end{bmatrix}$$
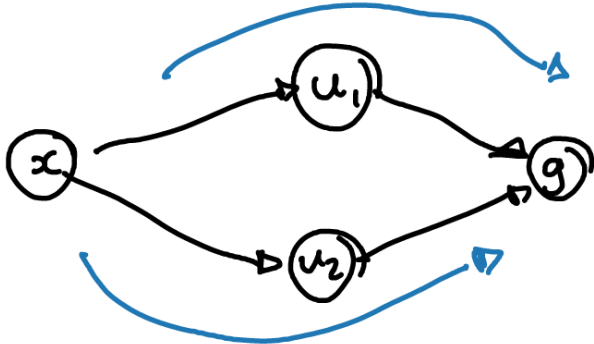
Using the above definitions, we can generalise the chain rule. Given $\mathbf{u}$ is a function of $\mathbf{x}$, and $g$ in turn is a vector function of $\mathbf{u}$, the vector-by-vector chain rule states:

$$\frac{\partial g(\mathbf{u})}{\partial \mathbf{x}} = \frac{\partial \mathbf{u}}{\partial \mathbf{x}} \frac{\partial g(\mathbf{u})}{\partial \mathbf{u}}$$

Order matters!

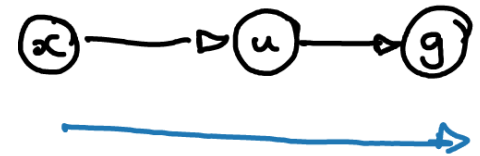$$\frac{\partial g_1}{\partial x_1} \begin{bmatrix} \circ \\ \end{bmatrix}$$

This generalised chain rule comes from the chain rule for multivariate functions. For scalars where $g$ depends on $u_1$ and $u_2$, which in turn depends on $x$, we have:

$$\frac{\partial g}{\partial x} = \frac{\partial u_1}{\partial x} \frac{\partial g}{\partial u_1} + \frac{\partial u_2}{\partial x} \frac{\partial g}{\partial u_2}$$
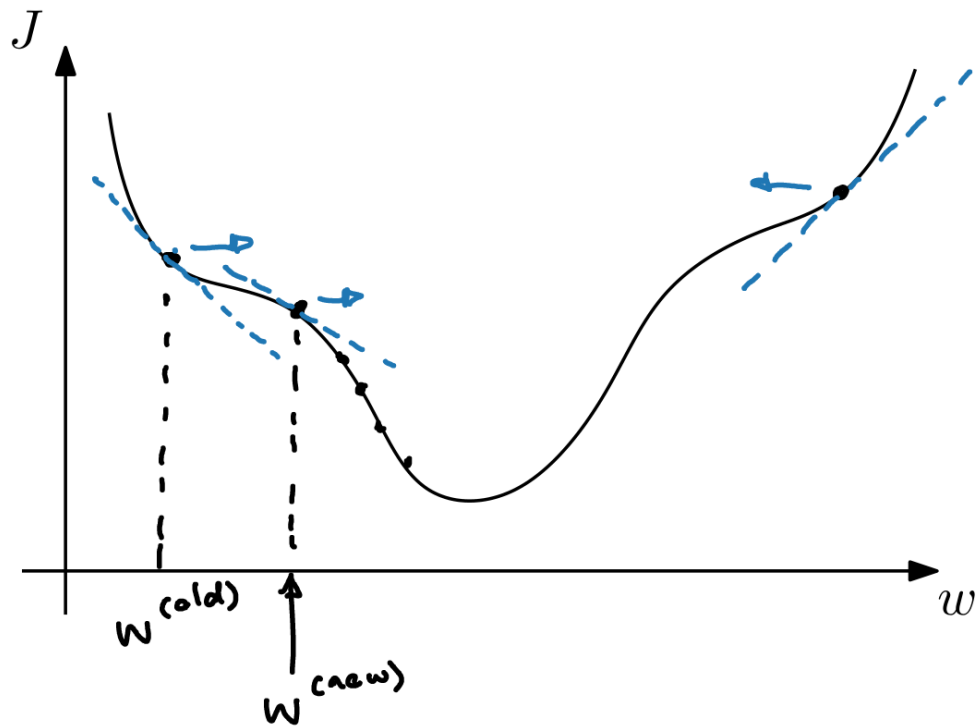
The chain rule of scalars:

$$\frac{dg}{dx} = \frac{du}{dx} \cdot \frac{dg}{du}$$

# Gradient descent recap



$$w \leftarrow w - \eta \frac{dJ}{dw}$$

$\uparrow$ Learning rate

$$\underline{w} \leftarrow \underline{w} - \eta \frac{\partial J}{\partial \underline{w}}$$

$$\underline{w}^{(new)} = \underline{w}^{(old)} - \eta \left. \frac{\partial J}{\partial \underline{w}} \right|_{\underline{w} = \underline{w}^{(old)}}$$
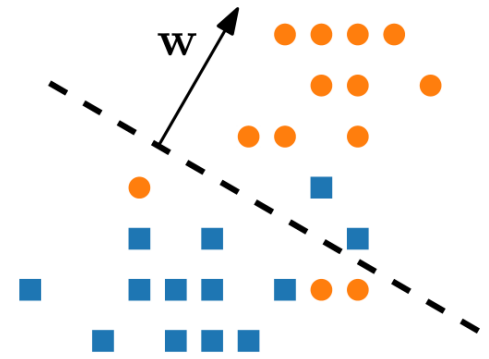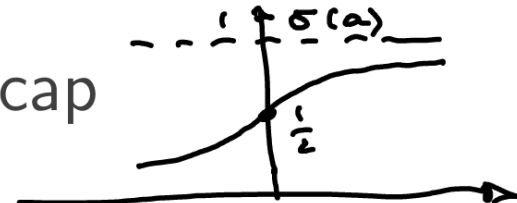
# Logistic regression, softmax regression, basis functions recap

- Prediction function: $f_{\mathbf{w}}(\mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x}) = \dfrac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}}$

- Interpret function as: $f_{\mathbf{w}}(\mathbf{x}) = P_{\mathbf{w}}(y = 1 | \mathbf{x})$

- With labels $y \in \{0, 1\}$, minimise the negative log likelihood:

$$(NLL)$$

$$J(\mathbf{w}) = -\log \prod_{n=1}^{N} P_{\mathbf{w}}(y^{(n)} | \mathbf{x}^{(n)})$$

$$= -\sum_{n=1}^{N} \left[ y^{(n)} \log f_{\mathbf{w}}(\mathbf{x}^{(n)}) + (1 - y^{(n)}) \log \left(1 - f_{\mathbf{w}}(\mathbf{x}^{(n)})\right) \right]$$

- Gradient: $\dfrac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = -\sum_{n=1}^{N} \left( y^{(n)} - f_{\mathbf{w}}(\mathbf{x}^{(n)}) \right) \mathbf{x}^{(n)}$

Softmax regression: $y \in \{1, 2, \ldots, K\}$

Basis functions:
$f_{\underline{w}}(\underline{x}) = \sigma(\underline{w}^\top \underline{\phi}(\underline{x}))$

# Introduction to neural networks

From logistic regression to neural networks

Herman Kamper

`http://www.kamperh.com/`
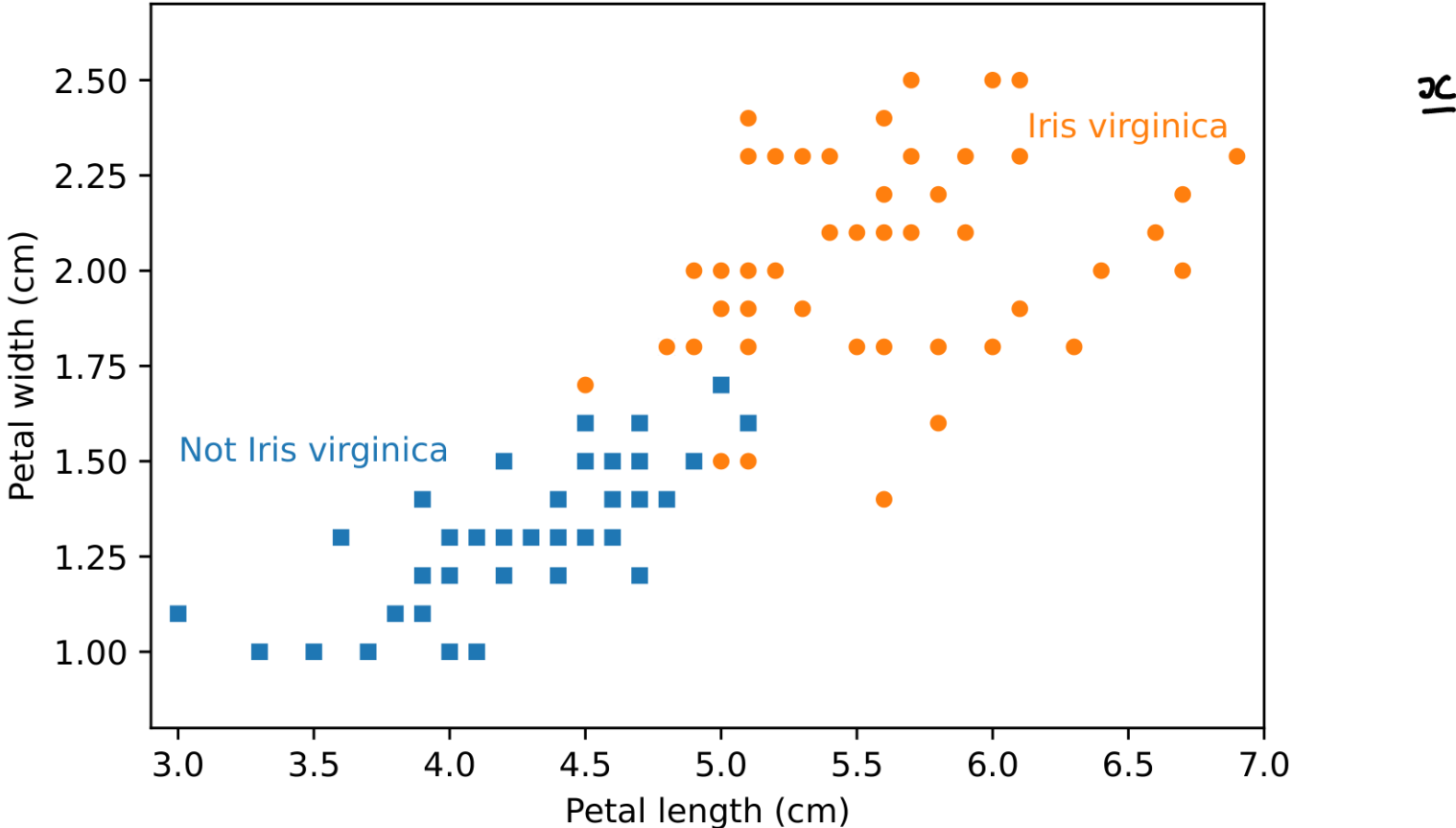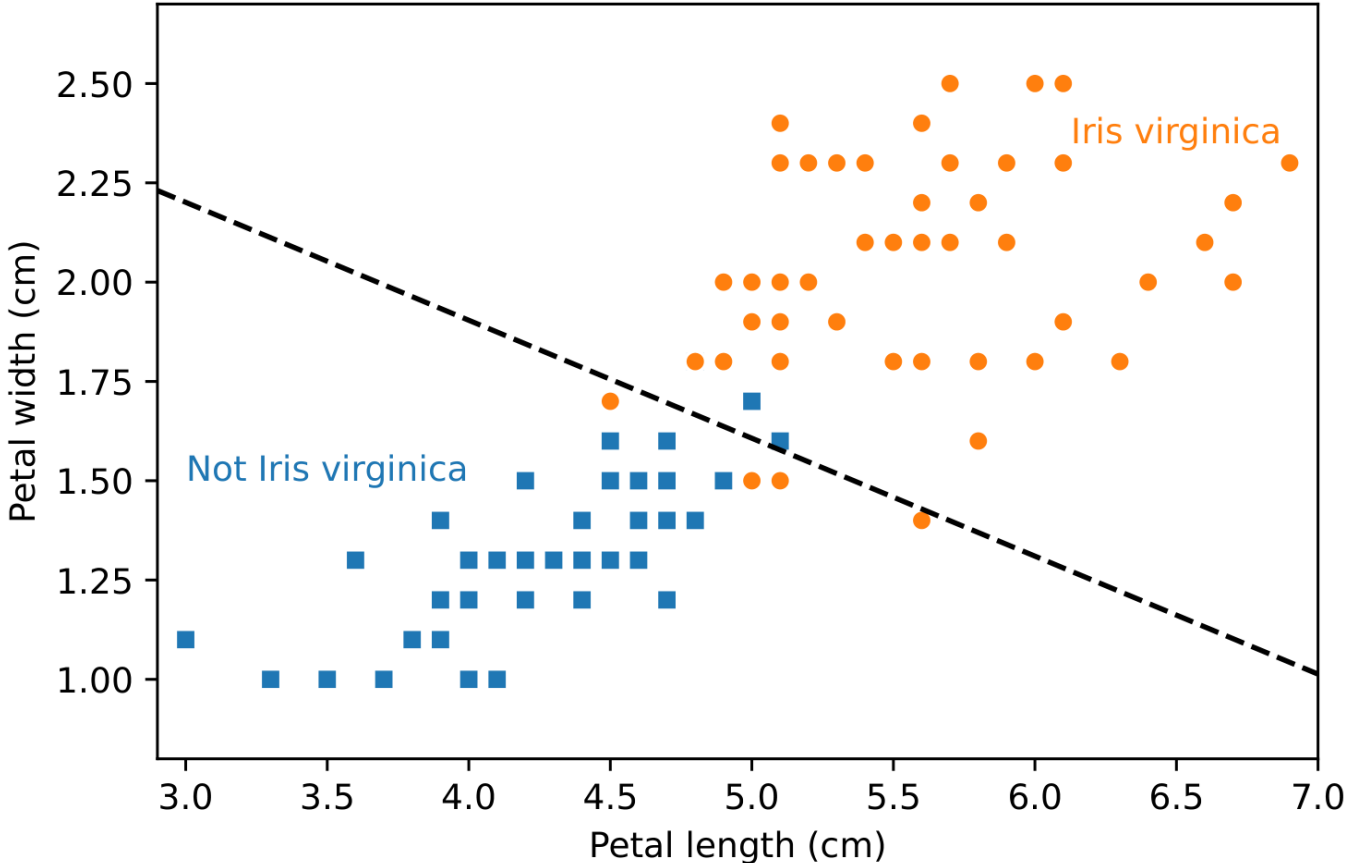
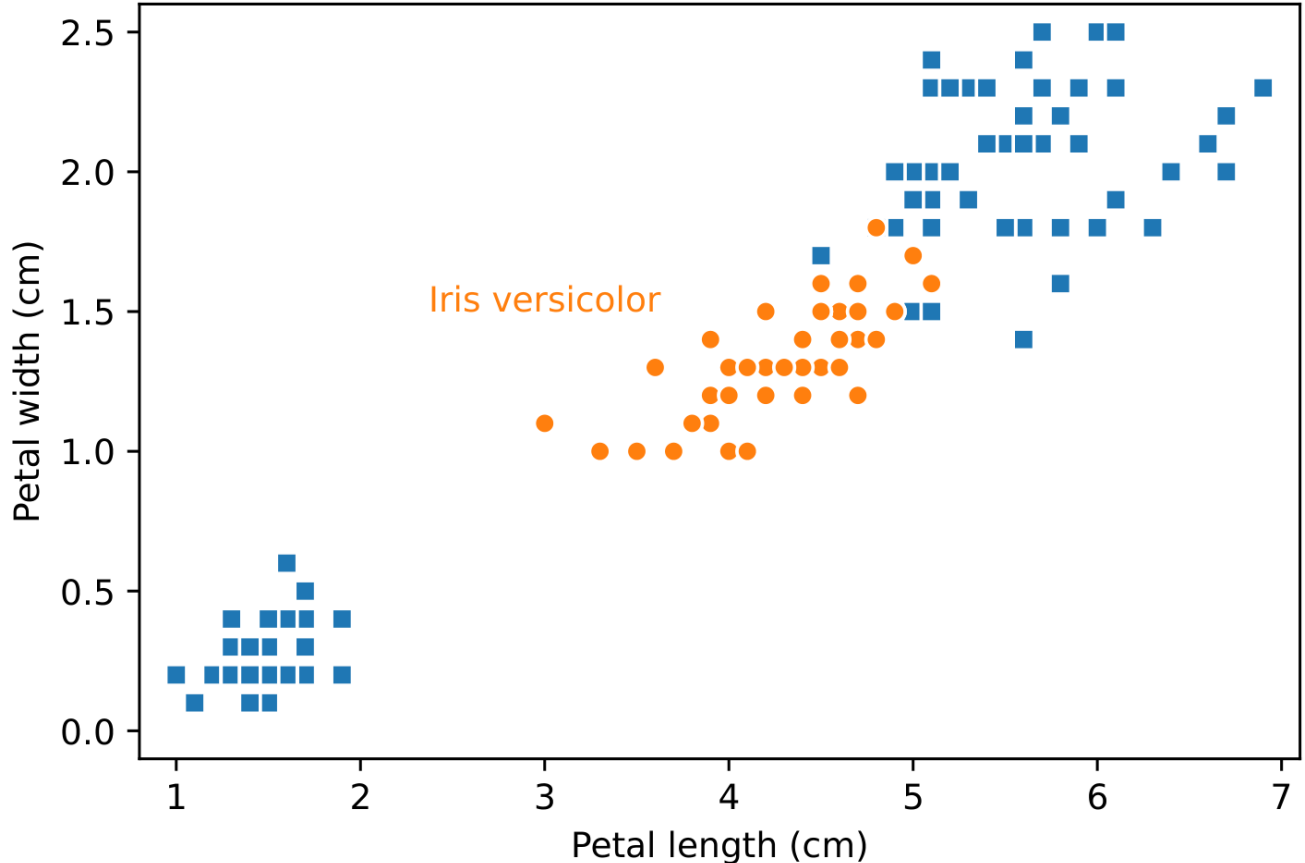# Binary classification of irises

# Logistic regression

# Logistic regression

# Nonlinear logistic regression

# Nonlinear logistic regression



$\phi(x)$

$x$

# Nonlinear logistic regression



$$\underline{\phi}(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1^2 \\ x_2^2 \end{bmatrix}$$

$$\underline{f}_w(x) = \sigma\left(\underline{w}^T \underline{\phi}(x)\right)$$

# Binary logistic regression with basis functions as a neural network



Logistic regression with basis functions:

$$f_{\underline{w}}(\underline{x}) = \sigma(\underline{w}^T \underline{\phi}(\underline{x})) = \frac{1}{1 + \exp(-\underline{w}^T \underline{\phi}(\underline{x}))}$$

We can set:

$$\phi_k(\underline{x}) = \sigma(\underline{w}_k^T \underline{x} + b_k)$$

or more generally:

$$\phi_k(\underline{x}) = g(\underline{w}_k^T \underline{x} + b_k)$$

Where $g$ is some non-linear function.

sigmoid          tanh          ReLU

$$f_{\underline{\theta}}(\underline{x}) = \hat{y} = g(\underline{w}^{[2]T} \underline{a}^{[1]} + b^{[2]})$$

$$\underline{\phi}(\underline{x}) = \underline{a}^{[1]} = g(\underline{W}^{[1]} \underline{x} + \underline{b}^{[1]})$$

$$\underline{\theta} = \{\underline{W}^{[1]}, \underline{b}^{[1]}, \underline{w}^{[2]}, b^{[2]}\}$$

# Why is it called a neural network?

(artificial)

# THE PERCEPTRON: A PROBABILISTIC MODEL FOR INFORMATION STORAGE AND ORGANIZATION IN THE BRAIN [1]

## F. ROSENBLATT

*Cornell Aeronautical Laboratory*

If we are eventually to understand the capability of higher organisms for perceptual recognition, generalization, recall, and thinking, we must first have answers to three fundamental and the stored pattern. According to this hypothesis, if one understood the code or "wiring diagram" of the nervous system, one should, in principle, be able to discover exactly what an

# Introduction to neural networks

Backpropagation (without forks)

Herman Kamper

`http://www.kamperh.com/`

# Example: Binary classification with a feedforward neural network

- How do we fit the parameters of our binary classification model?
  As usual: Use gradient descent to minimise the negative log likelihood.

$$\hat{y}^{(n)} = f_{\underline{\theta}}(x^{(n)})$$

- If we have a single training item $(\mathbf{x}^{(n)}, y^{(n)})$:

$$\hat{y}^{(n)} \in [0, 1]$$

$$J(\boldsymbol{\theta}) = -\left[ y^{(n)} \log \hat{y}^{(n)} + (1 - y^{(n)}) \log(1 - \hat{y}^{(n)}) \right] = \begin{cases} -\log \hat{y}^{(n)} & \text{if } y^{(n)} = 1 \\ -\log(1 - \hat{y}^{(n)}) & \text{if } y^{(n)} = 0 \end{cases}$$

- Need $\frac{\partial J}{\partial \mathbf{u}}$, where $\mathbf{u}$ is each of the parameters:

$$\underline{\theta} = \left\{ \mathbf{W}^{[1]}, \mathbf{b}^{[1]}, \mathbf{w}^{[2]}, b^{[2]} \right\}$$

$$\frac{\partial J}{\partial \underline{w}^{[2]}} = ?$$

- The *backpropagation algorithm* gives a principled procedure to obtain these gradients:
  Apply the chain rule while reusing previously computed results.

# The backpropagation algorithm (without forks)

- Represent your neural network as a computational graph.

- **Forward pass:** Start at the inputs and calculate the output of each operation in the graph. Store these values.

- **Backward pass:** Start at the output of the graph and move backwards. For each operation:

  (a) Determine and calculate the derivative of the output variable w.r.t. each of the input variables to the operation.

  (b) For each input variable $\mathbf{u}$, set

  $$\boldsymbol{\delta}_{\mathbf{u}} = \frac{\partial J}{\partial \mathbf{u}} = \frac{\partial \mathbf{z}}{\partial \mathbf{u}} \frac{\partial J}{\partial \mathbf{z}}$$

  where $\mathbf{z}$ is the output of the operation taking $\mathbf{u}$ as input.

- Use the calculated derivates to do take a gradient step to update the parameters. Repeat from the forward pass.

$$\frac{\partial z}{\partial v} \quad \text{and} \quad \frac{\partial z}{\partial u}$$

$$\frac{\partial J}{\partial \mathbf{u}} = \frac{\partial \mathbf{z}}{\partial \mathbf{u}} \frac{\partial J}{\partial \mathbf{z}}$$

## Example

Single training item:
$(\underline{x}^{(n)}, y^{(n)})$



(a) $-[y \log \hat{y} + (1-y) \log(1-\hat{y})] \leftarrow y$

(b) $\hat{y} = g(z^{[2]})$

$\sigma$    $g(\cdot)$

$z^{[2]} = \mathbf{w}^{[2]\top} \mathbf{a}^{[1]} + b^{[2]}$

(c) $\mathbf{w}^{[2]} \rightarrow$, $b^{[2]} \rightarrow$

$\mathbf{a}^{[1]}$

(d) ReLU    $g(\cdot)$

$z^{[1]} = \mathbf{W}^{[1]}\mathbf{x} + \mathbf{b}^{[1]}$

(e) $\mathbf{W}^{[1]} \rightarrow$, $\mathbf{b}^{[1]} \rightarrow$

$\mathbf{x}$

$\dfrac{\partial \underline{z}^{[1]}}{\partial \underline{W}^{[1]}} = ?$

(a)
$$\frac{\partial J}{\partial \hat{y}} = -\frac{\partial}{\partial \hat{y}}\left[y \log \hat{y} + (1-y)\log(1-\hat{y})\right]$$
$$= -\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}}$$
$$\delta_{\hat{y}} = \frac{\partial J}{\partial \hat{y}}\bigg|_{\substack{x=x^{(n)} \\ y=y^{(n)}}} = -\frac{y^{(n)}}{\hat{y}^{(n)}} + \frac{1-y^{(n)}}{1-\hat{y}^{(n)}}$$

(b)
$$\frac{\partial \hat{y}}{\partial z^{[2]}} = g'(z^{[2]})$$
$$\delta_{z^{[2]}} = \frac{\partial J}{\partial z^{[2]}}\bigg|_{\substack{x=x^{(n)} \\ y=y^{(n)}}} = \frac{\partial \hat{y}}{\partial z^{[2]}} \cdot \frac{\partial J}{\partial \hat{y}}$$
$$= g'(z^{[2]}) \cdot \delta_{\hat{y}}$$

(c)
$$\frac{\partial z^{[2]}}{\partial \underline{w}^{[2]}} = \frac{\partial}{\partial \underline{w}^{[2]}}\left[\underline{w}^{[2]\top}\underline{a}^{[1]} + b^{[2]}\right]$$
$$= \underline{a}^{[1]}$$
$$\delta_{\underline{w}^{[2]}} = \frac{\partial J}{\partial \underline{w}^{[2]}} = \frac{\partial z^{[2]}}{\partial \underline{w}^{[2]}} \cdot \frac{\partial J}{\partial z^{[2]}} = \underline{a}^{[1]} \cdot \delta_{z^{[2]}}$$
$$\delta_{b^{[2]}} = \frac{\partial J}{\partial b^{[2]}} = \frac{\partial z^{[2]}}{\partial b^{[2]}} \cdot \frac{\partial J}{\partial z^{[2]}}$$
$$\delta_{\underline{a}^{[1]}} = \frac{\partial J}{\partial \underline{a}^{[1]}} = \frac{\partial z^{[2]}}{\partial \underline{a}^{[1]}} \cdot \frac{\partial J}{\partial z^{[2]}}$$

(e) $\quad \underline{\delta}_{\underline{W}^{[i]}} = \dfrac{\partial J}{\partial \underline{W}^{[i]}} \overset{``}{=} \dfrac{\partial \underline{z}^{[i]}}{\partial \underline{W}^{[i]}} \cdot \dfrac{\partial J}{\partial \underline{z}^{[i]}} \overset{``}{}$

$$= \underline{\delta}_{\underline{z}^{[i]}} \, \underline{x}^{\intercal}$$

$\underline{\delta}_{\underline{b}^{[i]}} = \dfrac{\partial J}{\partial \underline{b}^{[i]}} = \dfrac{\partial \underline{z}^{[i]}}{\partial \underline{b}^{[i]}} \cdot \underline{\delta}_{\underline{z}^{[i]}}$

Why does this work?

$\underline{\delta}_{\underline{b}^{[i]}} = \dfrac{\partial J}{\partial \underline{b}^{[i]}} = \dfrac{\partial \underline{z}^{[i]}}{\partial \underline{b}^{[i]}} \cdot \dfrac{\partial J}{\partial \underline{z}^{[i]}} = \dfrac{\partial \underline{z}^{[i]}}{\partial \underline{b}^{[i]}} \cdot \dfrac{\partial \underline{a}^{[i]}}{\partial \underline{z}^{[i]}} \cdot \dfrac{\partial J}{\partial \underline{a}^{[i]}}$

$$= \dfrac{\partial \underline{z}^{[i]}}{\partial \underline{b}^{[i]}} \cdot \dfrac{\partial \underline{a}^{[i]}}{\partial \underline{z}^{[i]}} \cdot \dfrac{\partial \underline{z}^{[2]}}{\partial \underline{a}^{[i]}} \cdot \dfrac{\partial \hat{y}}{\partial \underline{z}^{[2]}} \cdot \dfrac{\partial J}{\partial \hat{y}}$$

Why the bar?

$\underline{\delta}_{\hat{y}} = \dfrac{\partial J}{\partial \hat{y}} \Bigg|_{\substack{x = x^{(n)} \\ y = y^{(n)} \\ \underline{\theta} = \hat{\underline{\theta}}^{(m)}}}$

$\hat{\underline{\theta}}^{(m+1)} = \hat{\underline{\theta}}^{(m)} - \gamma \dfrac{\partial J}{\partial \underline{\theta}} \Bigg|_{\underline{\theta} = \hat{\underline{\theta}}^{(m)}}$

# Binary classification of irises using our neural network

# Binary classification of irises using our neural network

# Multilayer feedforward neural network



$$\boldsymbol{\delta}_{\mathbf{z}^{[l]}} = \boldsymbol{\delta}_{\mathbf{a}^{[l]}} \odot g'(\mathbf{z}^{[l]})$$

$$\boldsymbol{\delta}_{\mathbf{b}^{[l]}} = \boldsymbol{\delta}_{\mathbf{z}^{[l]}}$$

$$\boldsymbol{\delta}_{\mathbf{W}^{[l]}} = \boldsymbol{\delta}_{\mathbf{z}^{[l]}} \mathbf{a}^{[l-1]\top}$$

$$\boldsymbol{\delta}_{\mathbf{a}^{[l-1]}} = \mathbf{W}^{[l]\top} \boldsymbol{\delta}_{\mathbf{z}^{[l]}}$$

These equations are sometimes combined:

$$\boldsymbol{\delta}_{\mathbf{z}^{[l]}} = \left( \mathbf{W}^{[l+1]\top} \boldsymbol{\delta}_{\mathbf{z}^{[l+1]}} \right) \odot g'(\mathbf{z}^{[l]})$$

# Introduction to neural networks

Computational graphs and automatic differentiation

Herman Kamper

http://www.kamperh.com/

# Why this (cool) graph formulation?

Adding additional structure is easy:

- As long as we know the derivative of a single operation, the gradient computation is fully specified by the graph.

- Each node just needs to know how to compute its output and how to compute the gradient w.r.t. its inputs ~~the gradient~~ ~~w.r.t. its output~~.

$\delta_x$   $x$

$z$

$\delta_z$

$\delta_y$   $y$

$\delta_x = \frac{\partial J}{\partial x}$

```python
class MultiplyGate():

    def forward(x, y):
        z = x*y
        self.x = x
        self.y = y
        return z


    def backward(delta_z):
        delta_x = self.y * delta_z
                # dz/dx  * dJ/dz
        delta_y = self.x * delta_z
                # dz/dy  * dJ/dz
        return [delta_x, delta_y]
```

# Why then study backprop if the software can do it?

- In some very simple cases, you might not want to have to rely on (the bulky) PyTorch or Tensorflow.
  E.g. the gradients for word2vec is relatively straightforward.

- Sometimes you might want to introduce a new computational operation and then you might need to implement the gradient computation.

- More often: You are hacking parts of the gradient computation for an existing block and need to modify it.

# Introduction to neural networks

Backpropagation (now general)

Herman Kamper

http://www.kamperh.com/

# A summary of derivatives for common blocks



$$\text{Scalars:} \quad \delta_v \quad \delta_z = \frac{\partial J}{\partial z}$$

$$\delta_u = \frac{\partial J}{\partial u}$$
$$= \frac{\partial z}{\partial u} \cdot \frac{\partial J}{\partial z}$$
$$= \frac{\partial z}{\partial u} \cdot \delta_z$$

$$\delta_b = \frac{\partial \mathbf{z}}{\partial \mathbf{b}} \, \delta_\mathbf{z}$$
$$= \mathbf{I} \, \delta_\mathbf{z} = \delta_\mathbf{z}$$

$$\delta_\mathbf{W} = \delta_\mathbf{z} \, \mathbf{x}^\top$$

$$\delta_\mathbf{x} = \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \, \delta_\mathbf{z}$$
$$= \mathbf{W}^\top \, \delta_\mathbf{z}$$

$$\delta_\mathbf{x} = \text{prod}\left(\frac{\partial \mathbf{z}}{\partial \mathbf{x}}, \delta_\mathbf{z}\right)$$

$$\delta_\mathbf{B} = \delta_\mathbf{C}^\top \mathbf{A}$$

$$\delta_\mathbf{B} = \text{prod}\left(\frac{\partial \mathbf{C}}{\partial \mathbf{B}}, \delta_\mathbf{C}\right)$$

$$\delta_\mathbf{A} = \delta_\mathbf{C} \mathbf{B}$$

# A general notation

In all of the above, for arbitrary variable U going in to operation with output Z, the error signal $\delta_U$ is obtained as some kind of product between $\frac{\partial Z}{\partial U}$ and $\frac{\partial J}{\partial Z} = \delta_Z$.

- With vectors as inputs and outputs: $\delta_\mathbf{b} = \frac{\partial \mathbf{z}}{\partial \mathbf{b}} \delta_\mathbf{z}$

- But sometimes the order between $\frac{\partial Z}{\partial U}$ and $\delta_Z$ flips: $\delta_\mathbf{W} = \delta_\mathbf{z} \mathbf{x}^\top$

- Or we have to take the transpose: $\delta_\mathbf{B} = \delta_\mathbf{C}^\top \mathbf{A}$

Let's capture all of these with the new prod operator:

$$\delta_U = \frac{\partial J}{\partial U} = \text{prod}\left(\frac{\partial Z}{\partial U}, \frac{\partial J}{\partial Z}\right)$$

$$= \text{prod}\left(\frac{\partial Z}{\partial U}, \delta_Z\right)$$

d2l.ai

$\delta_z$

$u \longrightarrow z$

$\delta_u$

$\delta_u = \frac{\partial J}{\partial u} = \frac{\partial z}{\partial u} \cdot \frac{\partial J}{\partial z}$

# About forks

$$\delta_u = \frac{\partial J}{\partial u} = \frac{\partial z}{\partial u} \cdot \frac{\partial J}{\partial z} + \frac{\partial a}{\partial u} \cdot \frac{\partial J}{\partial a}$$

$$= \frac{\partial z}{\partial u} \cdot \delta_z + \frac{\partial a}{\partial u} \cdot \delta_a$$

"Error signal" $\Big\}$
"Accumulator" $\Big\}$ → $\delta_u$



$\dfrac{\partial \mathbf{a}}{\partial \mathbf{u}} \dfrac{\partial J}{\partial \mathbf{a}}$

$\delta_a$

$\dfrac{\partial \mathbf{z}}{\partial \mathbf{u}} \dfrac{\partial J}{\partial \mathbf{z}}$

$\delta_z$

In general:

$$\delta_u = \mathrm{prod}\left(\frac{\partial z}{\partial u}, \delta_z\right) + \mathrm{prod}\left(\frac{\partial A}{\partial u}, \delta_A\right)$$

Example: $L_2$ regularisation

$$J(\underline{\theta}) = -\left[y^{(n)} \log \hat{y}^{(n)} + (1 - y^{(n)}) \cdot \log (1 - \hat{y}^{(n)})\right]$$
$$+ \lambda \underline{w}^{[2]T} \underline{w}^{[2]}$$

$$\underline{\delta}_{\underline{w}^{[2]}} = \frac{\partial J}{\partial \underline{w}^{[2]}} = \frac{\partial z^{[2]}}{\partial \underline{w}^{[2]}} \cdot \delta_{z^{[2]}} + \frac{\partial r}{\partial \underline{w}^{[2]}} \cdot \delta_r$$

$$= \underline{a}^{[\cdot]} \delta_{z^{[2]}} + 2\lambda \underline{w}^{[2]} \delta_r$$

$\delta_r$

$J$

$+$

$r = \lambda \mathbf{w}^{[2]\top} \mathbf{w}^{[2]}$

$L_2$    $r$    $-[y \log \hat{y} + (1 - y) \log (1 - \hat{y})]$

$\hat{y}$

$g(\cdot)$

$\delta_{z^{[2]}}$

$z^{[2]} = \mathbf{w}^{[2]\top} \mathbf{a}^{[1]} + b^{[2]}$

$\mathbf{w}^{[2]}$
$b^{[2]}$

$\mathbf{a}^{[1]}$

$g(\cdot)$

$\mathbf{z}^{[1]} = \mathbf{W}^{[1]} \mathbf{x} + \mathbf{b}^{[1]}$

$\mathbf{W}^{[1]}$
$\mathbf{b}^{[1]}$

$\mathbf{x}$

# The backpropagation algorithm (now with forks)

- **Initialisation:**
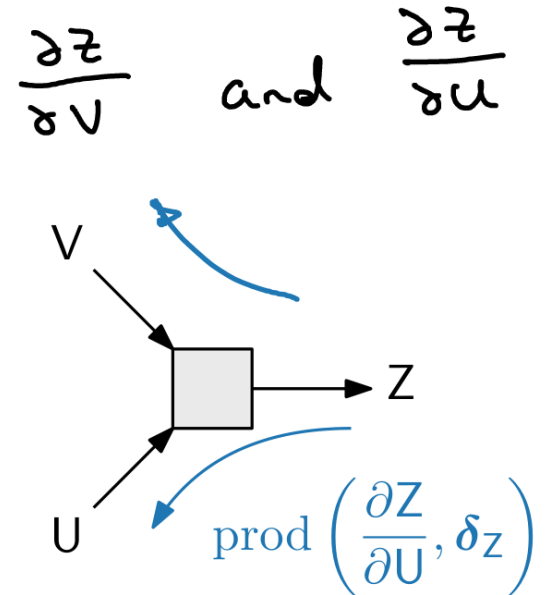  Set accumulators to zero for all input variables: $\delta_U \leftarrow \mathbf{0}$

- **Forward pass:** Start at the inputs and calculate the output of each operation in the graph. Store these values.

- **Backward pass:** Start at the output of the graph and move backwards. For each operation:

  (a) Determine and calculate the derivative of the output variable w.r.t. each of the input variables to the operation.

  (b) For each input variable U, add $\mathrm{prod}\left(\dfrac{\partial Z}{\partial U}, \dfrac{\partial J}{\partial Z}\right)$ to its accumulator, i.e.

  $$\delta_U \leftarrow \delta_U + \mathrm{prod}\left(\frac{\partial Z}{\partial U}, \delta_Z\right)$$

$$\frac{\partial z}{\partial v} \quad \text{and} \quad \frac{\partial z}{\partial u}$$

V

Z

U $\qquad$ $\mathrm{prod}\left(\dfrac{\partial Z}{\partial U}, \delta_Z\right)$

# Introduction to neural networks

Relationship between negative log likelihood and cross entropy

Herman Kamper

http://www.kamperh.com/

# Negative log likelihood and cross entropy

$$\hat{y} = \underline{f}_{\underline{\theta}}(\underline{x}) = \text{softmax}(\underline{z}) = \frac{1}{\sum_{j=1}^{k} e^{z_j}} \begin{bmatrix} e^{z_1} \\ e^{z_2} \\ \vdots \\ e^{z_k} \end{bmatrix}$$

$$y^{(n)} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} k \Bigg\} K$$

$\underline{f}_{\underline{\theta}}(\underline{x})$ [o o o]
↑ softmax
[o o o]  $\underline{z}$

$\underline{x}$

NLL: 
$$J(\underline{\theta}) = -\log \prod_{n=1}^{N} P_{\underline{\theta}}(y^{(n)} | \underline{x}^{(n)})$$

$$= -\sum_{n=1}^{N} \log P_{\underline{\theta}}(y^{(n)} | \underline{x}^{(n)})$$

$$= -\sum_{n=1}^{N} \sum_{k=1}^{K} y_k^{(n)} \log \frac{e^{z_k}}{\sum_{j=1}^{K} e^{z_j}}$$

$$= -\sum_{n=1}^{N} \sum_{k=1}^{K} y_k^{(n)} \log \hat{y}_k^{(n)}$$

$$\boxed{\text{Cross-entropy}: H(p,q) = -\sum_{k=1}^{K} p_k \log q_k}$$

# Introduction to neural networks

Neural networks in practice and NLP examples

Herman Kamper

http://www.kamperh.com/

# The art of neural networks

- Sometimes useful to scale inputs.

- Instead of vanilla gradient descent, we often use advanced forms of mini-batch gradient descent (Adam is popular at the moment).

- Different initialisation strategies, e.g. `https://arxiv.org/abs/1811.00293`.

- Overfitting: Can combat using standard regularisation, but often rather just use dropout or rely on SGD with early stopping.

- Need to choose number of hidden layers and number of units per layer, and often many more hyperparameters.

- Often make architecture choices (e.g. skip connections) to deal with optimisation problems (e.g. exploding or vanishing gradients).

# Named entity recognition

```
last night Paris  Hilton wowed in a sequin gown
              PER    PER

Samuel Quinn was arrested in the Hilton Hotel  in Paris in April 1989
PER    PER                       LOC    LOC       LOC       DATE  DATE
```

| Tag  | Description                    | Example                                        |
| ---- | ----------------------------- | ---------------------------------------------- |
| PER  | People, characters            | **Shannon** is a giant of information theory.  |
| ORG  | Organisation                  | The **ICC** is the governing body of cricket.  |
| LOC  | Location                      | **Mt. Sanitas** is in **Sunshine Canyon**.     |
| GPE  | Geo-political (countries, states) | Petrol prices are going up in **South Africa**. |
| DATE | Days, months, years           | Micah was born in **April**.                   |

We want to classify the entity of `Paris` in the sentence:
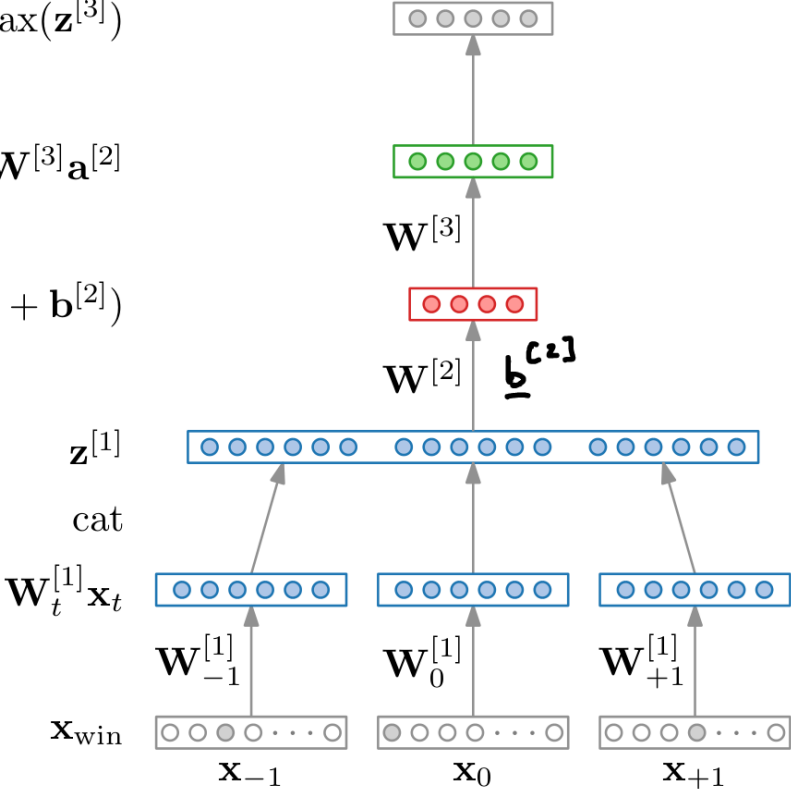
`anywhere in Paris museums are great`

Use a window of words around the centre word:

$$\mathbf{x}_{\mathrm{win}} = \begin{bmatrix} -\mathbf{x}_{\mathrm{in}}^{\top}- & -\mathbf{x}_{\mathrm{Paris}}^{\top}- & -\mathbf{x}_{\mathrm{museums}}^{\top}- \end{bmatrix}^{\top}$$

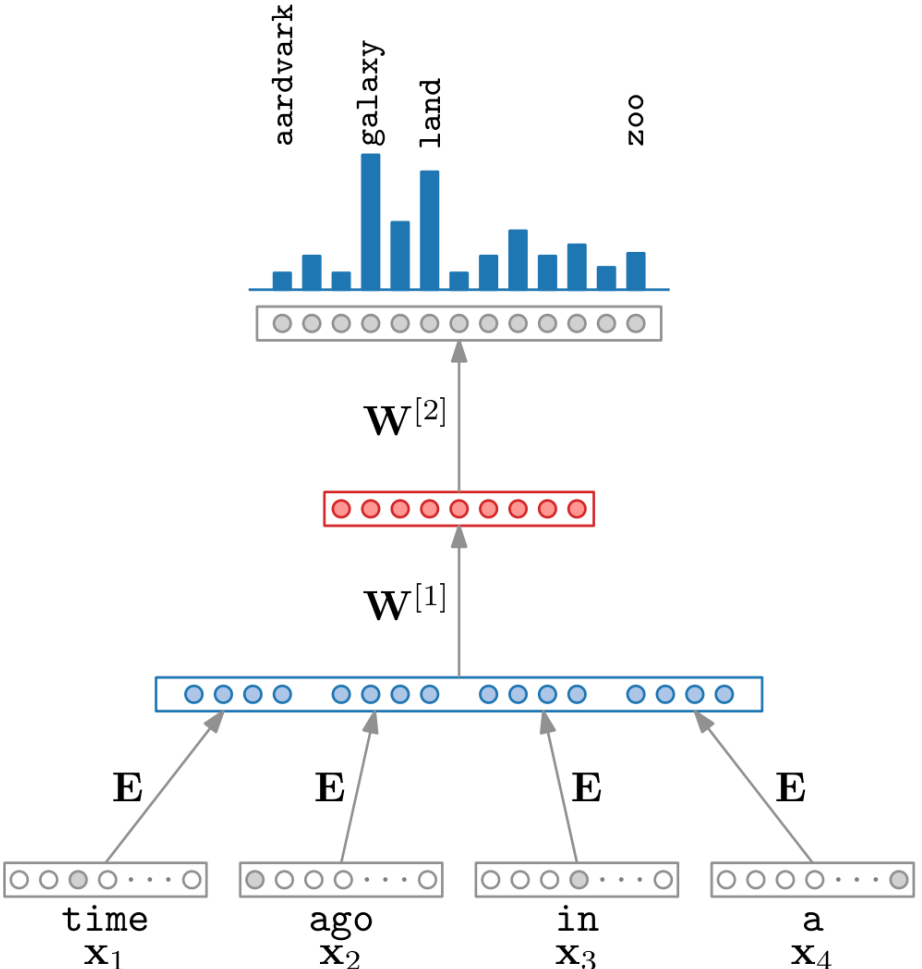$$f_{\boldsymbol{\theta}}(\mathbf{x}_{\mathrm{win}}) = \mathrm{softmax}(\mathbf{z}^{[3]})$$

$$\mathbf{z}^{[3]} = \mathbf{W}^{[3]}\mathbf{a}^{[2]}$$

$$\mathbf{a}^{[2]} = g(\mathbf{W}^{[2]}\mathbf{z}^{[1]} + \mathbf{b}^{[2]})$$

$\mathbf{z}^{[1]}$

cat

$\mathbf{W}_{t}^{[1]}\mathbf{x}_{t}$

$\mathbf{W}^{[3]}$

$\mathbf{W}^{[2]} \quad \underline{b}^{[2]}$

$\mathbf{W}_{-1}^{[1]} \qquad \mathbf{W}_{0}^{[1]} \qquad \mathbf{W}_{+1}^{[1]}$

$\mathbf{x}_{\mathrm{win}}$

$\mathbf{x}_{-1} \qquad \mathbf{x}_{0} \qquad \mathbf{x}_{+1}$

# Neural language models

`A long long time ago in a ...`



Output distribution
$$f_{\boldsymbol{\theta}}(\mathbf{x}_{1:4}) = \hat{\mathbf{y}}$$
$$= \mathrm{softmax}(\mathbf{W}^{[2]}\mathbf{h} + \mathbf{b}^{[2]})$$
$$\in [0,1]^{|\mathcal{V}|}$$

Hidden layer
$$\mathbf{h} = g(\mathbf{W}^{[1]}\mathbf{e} + \mathbf{b}^{[1]})$$

Concatenated word embeddings
$$\mathbf{e} = [\mathbf{e}_1; \mathbf{e}_2; \mathbf{e}_3; \mathbf{e}_4]$$

One-hot word vectors

# Introduction to neural networks

### What did I read

Herman Kamper

`http://www.kamperh.com/`

- H. Kamper, "Yet another introduction to backpropagation," *Stellenbosch University*, 2022.
  http://www.kamperh.com/notes/kamper_backprop22.pdf

- A. Ng, "Multi-layer neural network," *Stanford University*, 2013.
  http://ufldl.stanford.edu/tutorial/supervised/MultiLayerNeuralNetworks/

- A. Zhang, Z. C. Lipton, M. Li, and A. J. Smola, *Dive into Deep Learning*, 2021.
  https://d2l.ai/

- M. P. Deisenroth, A. A. Faisal, and C. S. Ong, *Mathematics for Machine Learning*, 2020.
  https://mml-book.github.io/book/mml-book.pdf

- CS231n: Optimization 2, *Stanford University*.
  https://cs231n.github.io/optimization-2/

- C. Manning, "CS224N: Neural net learning, gradients by hand (matrix calculus) and algorithmically (the backpropagation algorithm)," *Stanford University*, 2022.
  https://web.stanford.edu/class/cs224n/