

Business Forecasting

ADIA Course

Day 4 – Session 1

Dr. Tanujit Chakraborty

Sorbonne University, Abu Dhabi, UAE

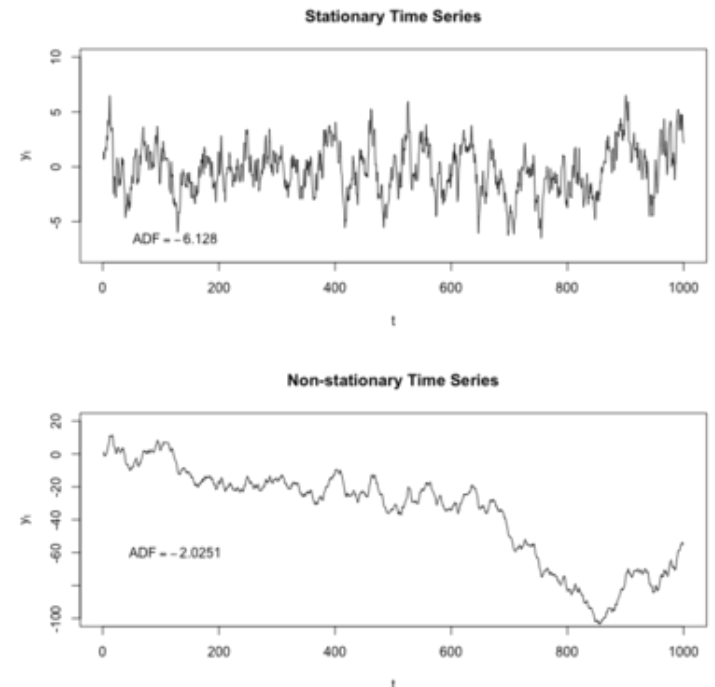
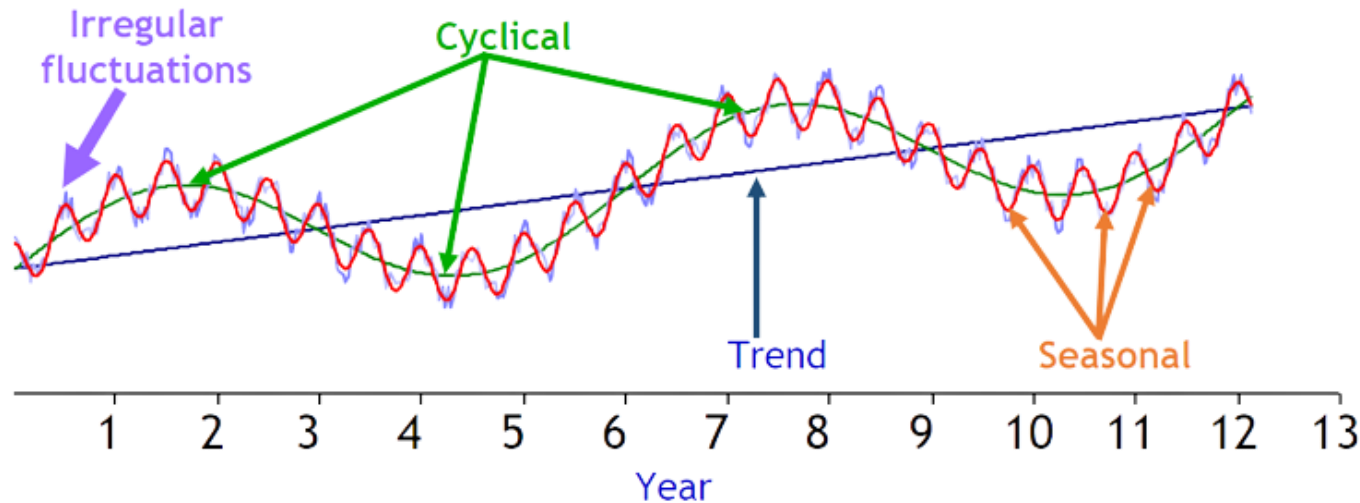
Sorbonne Centre for Artificial Intelligence, Paris, France

Research Areas: Time Series Forecasting, Machine Learning, Econometrics, Health Data Science

July 3, 2025

Learning from Time-Series Data

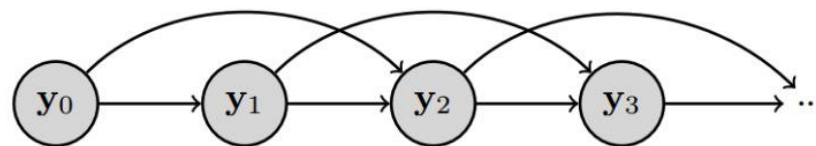
- The input is a sequence of (non-i.i.d.) examples y_1, y_2, \dots, y_t .
- The problem may be supervised or unsupervised, e.g.,
 - Forecasting: Predict y_{t+1} using y_1, y_2, \dots, y_t
 - Cluster the examples or perform dimensionality reduction / Anomaly detection
- Evolution of time-series data can be attributed to several factors



- Teasing apart these factors of variation is also an important problem.

- Auto-regressive (AR): Regress each example on p previous lagged values - AR(p) model

$$y_t = c + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \cdots + \phi_p y_{t-p} + \varepsilon_t,$$



Auto-regressive Model (shown above: 2nd order AR)

- Moving Average (MA): Regress each example on q previous stochastic errors - MA(q) model

$$y_t = c + \varepsilon_t + \theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + \cdots + \theta_q \varepsilon_{t-q},$$

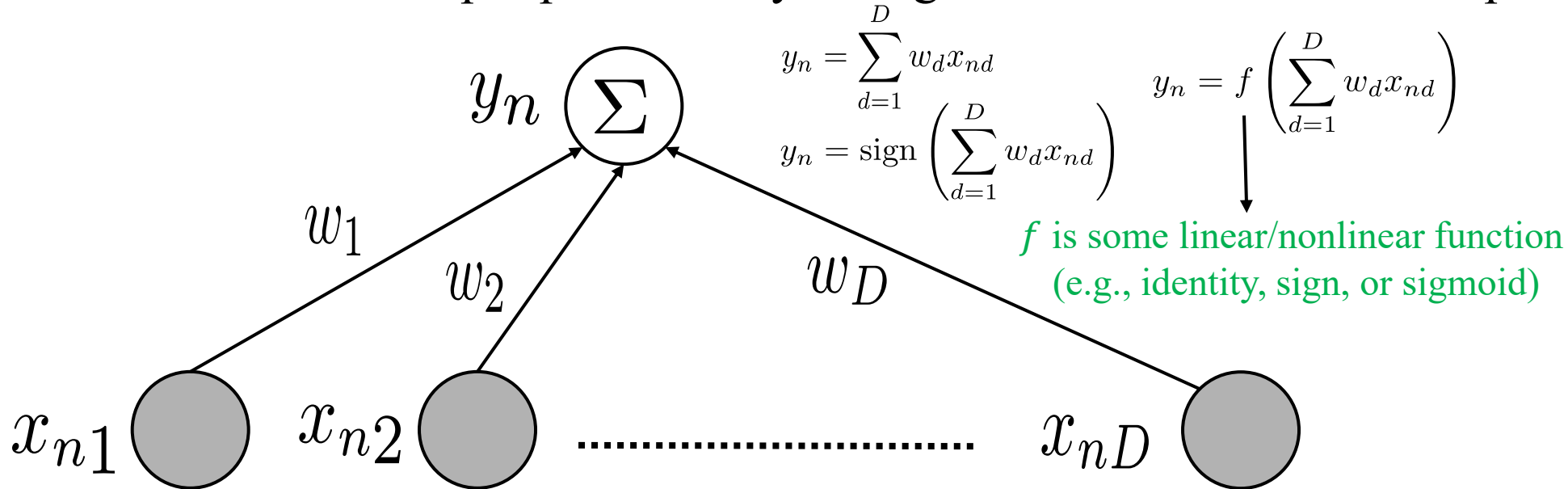
- Auto-regressive Integrated Moving Average (ARMA): Regress each example of p previous lagged values and q previous stochastic errors

$$y'_t = c + \phi_1 y'_{t-1} + \cdots + \phi_p y'_{t-p} + \theta_1 \varepsilon_{t-1} + \cdots + \theta_q \varepsilon_{t-q} + \varepsilon_t,$$

where y'_t is the differenced series (if the data is nonstationary and differencing is applied). We call this an **ARIMA**(p, d, q) model.

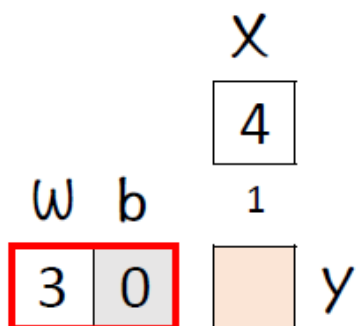
Limitation of Linear Models

- Linear models: Output produced by taking a linear combination of input features

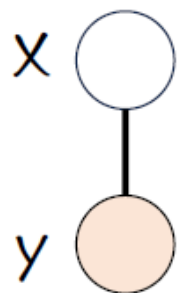


- A basic unit of the form $y = f(\mathbf{w}^\top \mathbf{x})$ is known as the “Perceptron” (not to be confused with the Perceptron “algorithm”, which learns a linear classification model)
- This can’t however learn nonlinear functions or nonlinear decision boundaries

Exercise: Linear Layer

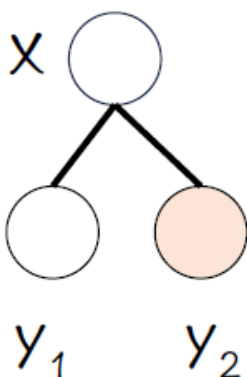
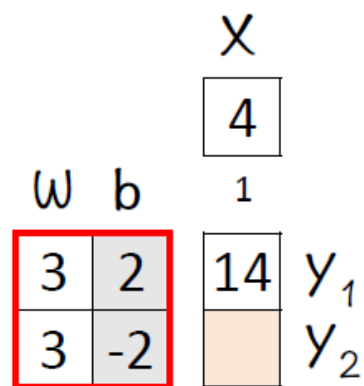


$$y = [w|b] \cdot [X|1] = wX + b$$



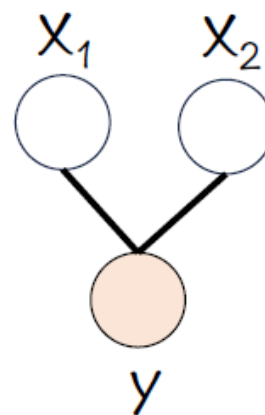
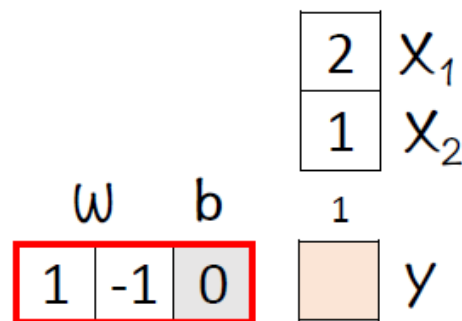
Exercise 1

12



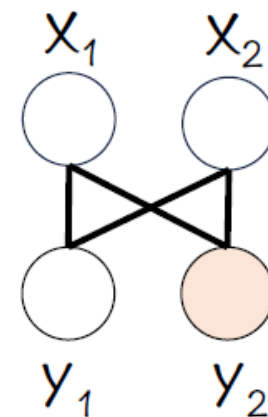
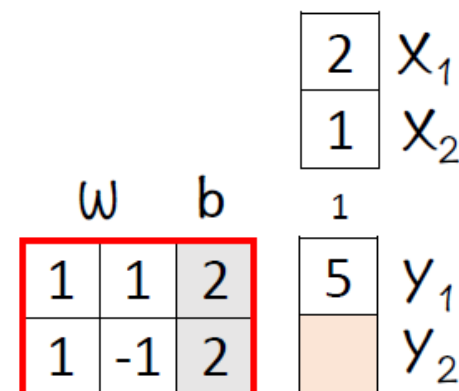
Exercise 2

10



Exercise 3

1

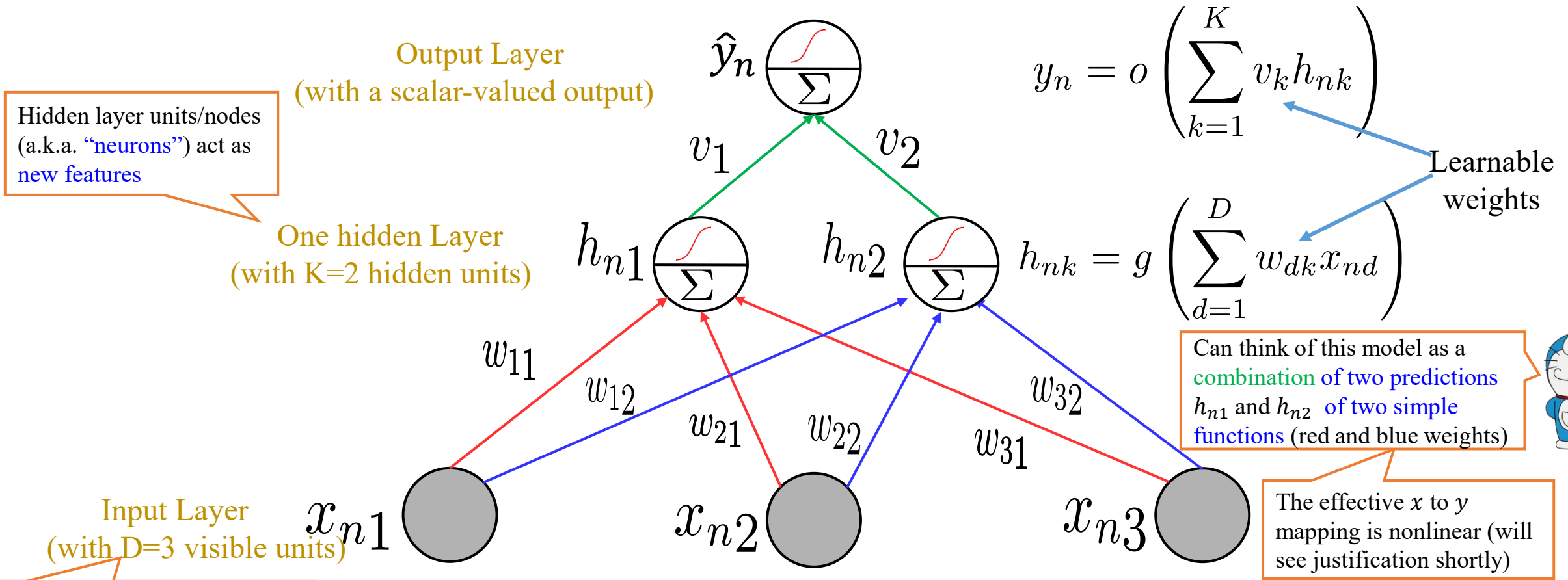


Exercise 4

3

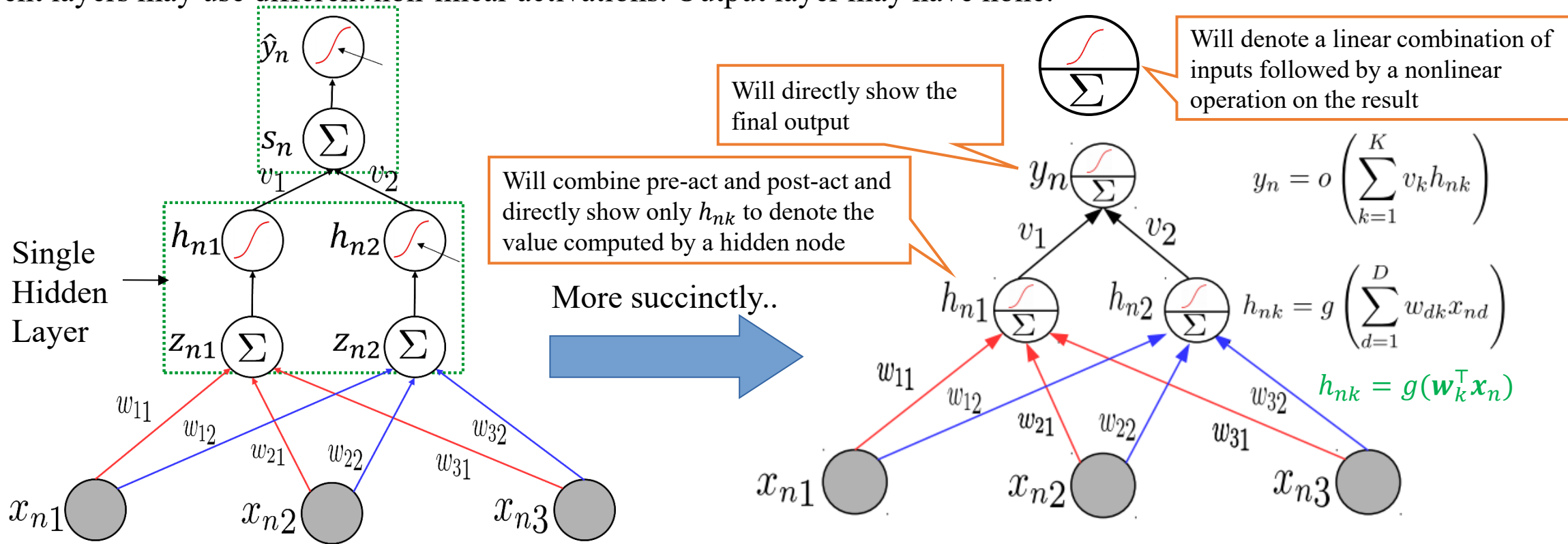
Neural Networks: Multi-layer Perceptron (MLP)

- An MLP consists of an **input layer**, an **output layer**, and **one or more hidden layers**



Neural Nets: A Compact Illustration

- Note: Hidden layer pre-act z_{nk} and post-act h_{nk} will be shown together for brevity
- Different layers may use different non-linear activations. Output layer may have none.

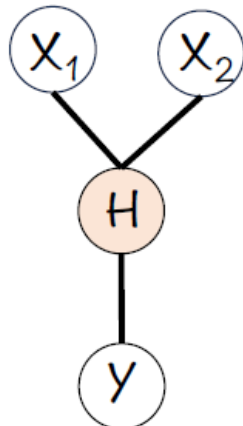


- Denoting $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K]$, $\mathbf{w}_k \in \mathbb{R}^D$, $\mathbf{h}_n = g(\mathbf{W}^\top \mathbf{x}_n) \in \mathbb{R}^K$ ($K = 2, D = 3$ above).
- **Note:** g applied elementwise on pre-activation vector $\mathbf{z}_n = \mathbf{W}^\top \mathbf{x}_n$

Exercise: Hidden Layer

$$\begin{matrix} w & b \\ \boxed{1} & \boxed{2} & \boxed{0} \end{matrix}$$

$$\begin{matrix} 1 \\ \boxed{1} \\ -1 \end{matrix} \begin{matrix} X_1 \\ X_2 \\ 1 \end{matrix} \xrightarrow{\text{ReLU}} \begin{matrix} \boxed{2} \\ -1 \end{matrix} \approx \begin{matrix} \boxed{0} \end{matrix} y$$

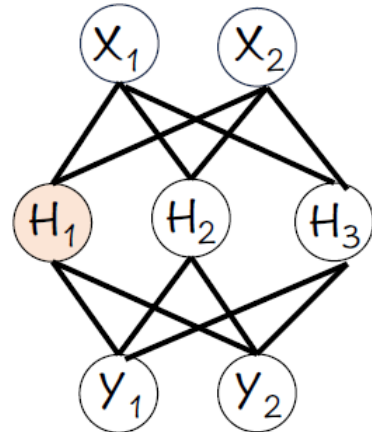


Exercise 1

$$\begin{matrix} \boxed{-1} & \boxed{0} \end{matrix}$$

$$\begin{matrix} 3 \\ 1 \end{matrix} \begin{matrix} X_1 \\ X_2 \end{matrix} \xrightarrow{\text{ReLU}} \begin{matrix} \boxed{0} & \boxed{2} & \boxed{1} \\ 1 & 1 & -1 \\ 1 & 0 & -1 \end{matrix} \approx \begin{matrix} \boxed{3} \\ 2 \end{matrix} \begin{matrix} H_1 \\ H_2 \\ H_3 \end{matrix}$$

$$\begin{matrix} \boxed{1} & \boxed{1} & \boxed{0} & \boxed{0} \\ 0 & 1 & 1 & -1 \end{matrix} \begin{matrix} 6 \\ 4 \end{matrix} \xrightarrow{\text{ReLU}} \begin{matrix} \boxed{6} \\ 4 \end{matrix} \begin{matrix} Y_1 \\ Y_2 \end{matrix}$$

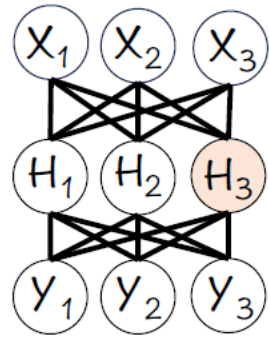


Exercise 2

$$\begin{matrix} \boxed{3} & \boxed{3} \end{matrix}$$

$$\begin{matrix} 1 \\ 0 \\ 2 \end{matrix} \begin{matrix} X_1 \\ X_2 \\ X_3 \end{matrix} \xrightarrow{\text{ReLU}} \begin{matrix} \boxed{0} & \boxed{1} & \boxed{1} & \boxed{-1} \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & -1 \end{matrix} \approx \begin{matrix} \boxed{1} \\ 3 \end{matrix} \begin{matrix} H_1 \\ H_2 \\ H_3 \end{matrix}$$

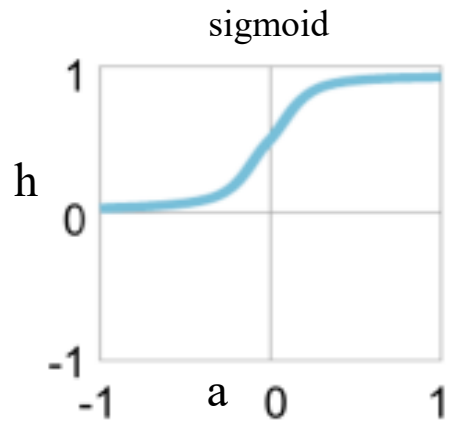
$$\begin{matrix} \boxed{1} & \boxed{1} & \boxed{0} & \boxed{0} \\ 0 & -1 & 0 & 1 \\ 1 & 1 & 1 & -1 \end{matrix} \begin{matrix} 4 \\ -2 \\ 3 \end{matrix} \xrightarrow{\text{ReLU}} \begin{matrix} \boxed{4} \\ 0 \\ 3 \end{matrix} \begin{matrix} Y_1 \\ Y_2 \\ Y_3 \end{matrix}$$



Exercise 3

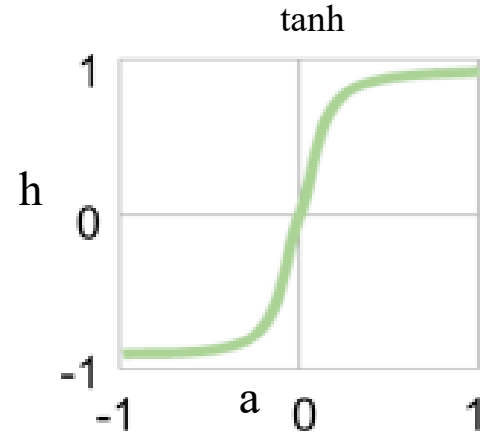
$$\begin{matrix} \boxed{-1} & \boxed{0} \end{matrix}$$

Activation Functions: Some Common Choices



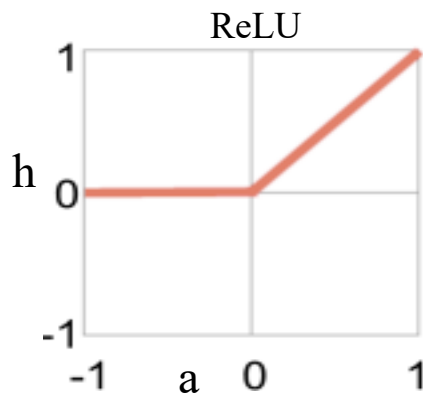
For sigmoid as well as tanh, gradients saturate (become close to zero as the function tends to its extreme values)

Sigmoid: $h = \sigma(a) = \frac{1}{1 + \exp(-a)}$



Preferred more than sigmoid. Helps keep the mean of the next layer's inputs close to zero (with sigmoid, it is close to 0.5)

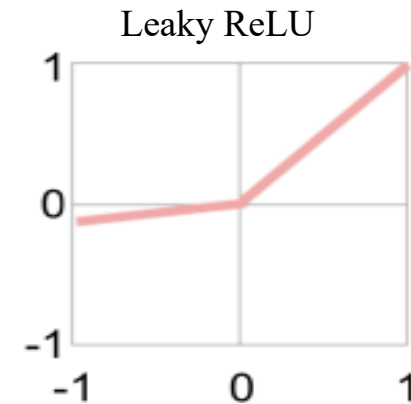
tanh (tan hyperbolic): $h = \frac{\exp(a) - \exp(-a)}{\exp(a) + \exp(-a)} = 2\sigma(2a) - 1$



ReLU and Leaky ReLU are among the most popular ones (also efficient to compute)

Helps fix the dead neuron problem of ReLU when a is a negative number

ReLU (Rectified Linear Unit): $h = \max(0, a)$

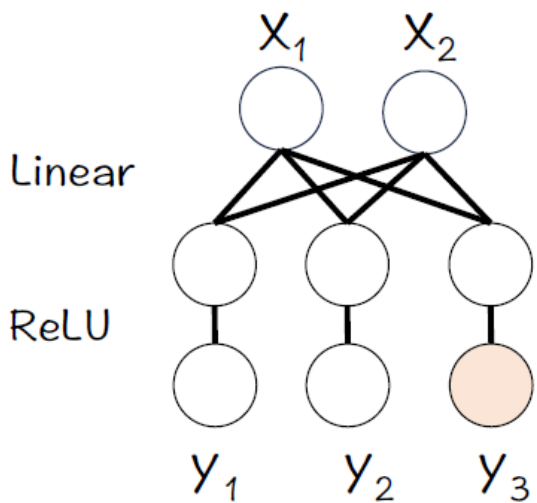


Without nonlinear activation, a deep neural network is equivalent to a linear model no matter how many layers we use

Leaky ReLU: $h = \max(\beta a, a)$
where β is a small positive number

Exercises: Activation Functions-ReLU, Sigmoid, Tanh

			3	X_1			
			2	X_2			
w	b		1				
1	-1	0	1	ReLU	1	y_1	
1	-2	0	-1		0	y_2	
1	-3	0	-3			y_3	



Exercise 1

0

			3	X_1			
			2	X_2			
w	b		1				
1	-1	0	1	Sigmoid $\sigma\{3\}$	0.5	y_1	
1	-2	0	-1			y_2	
1	-3	0	-3		0	y_3	

3 Level Quantization
for hand calculation

Integer	Sigmoid
≥ 2	1
-1, 0, 1	0.5
≤ -2	0

Exercise 2

0.5

			3	X_1			
			2	X_2			
w	b		1				
1	-1	-2	-1	Tanh{3}	0	y_1	
1	-2	-2	-3			y_2	
1	-3	-2	-5		-1	y_3	

3 Level Quantization
for hand calculation

Integer	Sigmoid	Tanh
≥ 2	1	1
-1, 0, 1	0.5	0
≤ -2	0	-1

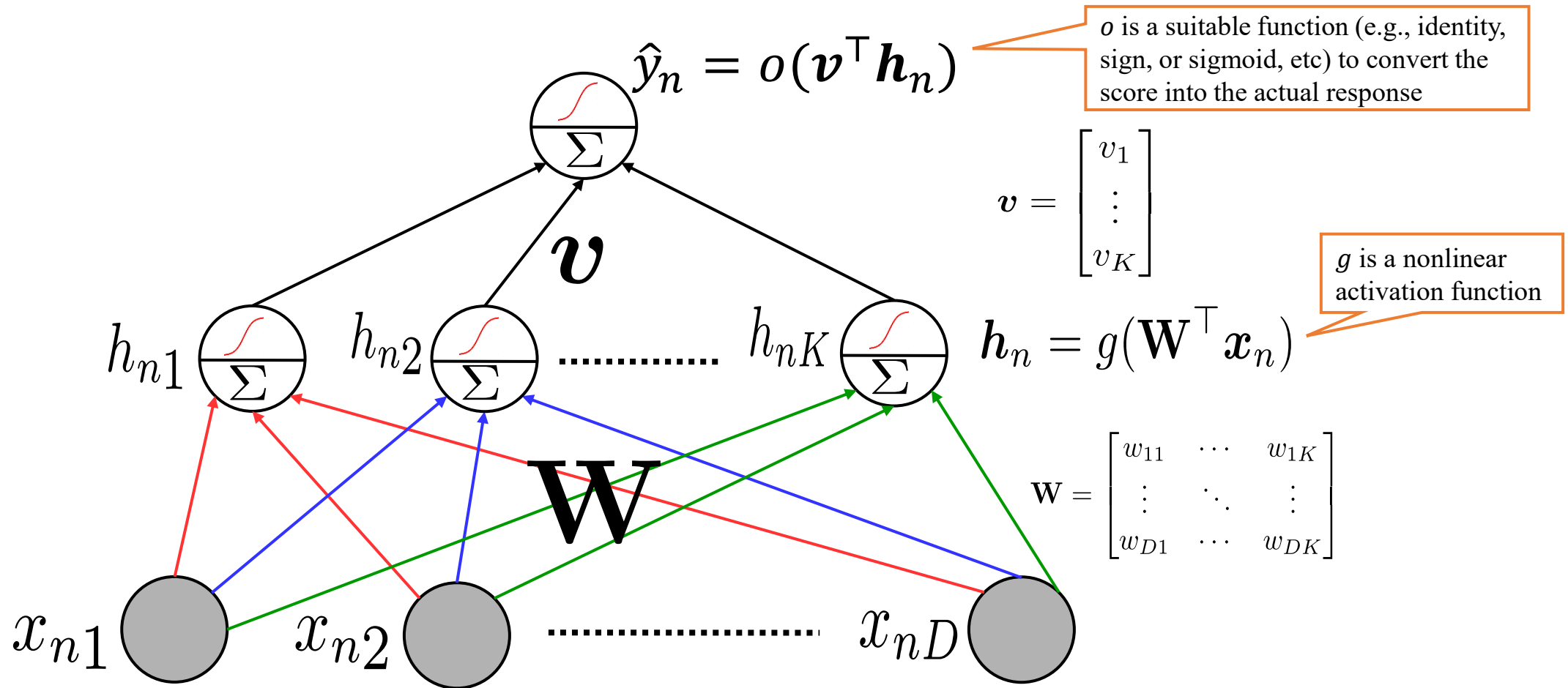
Exercise 3

-1

Examples of some basic NN/MLP architectures

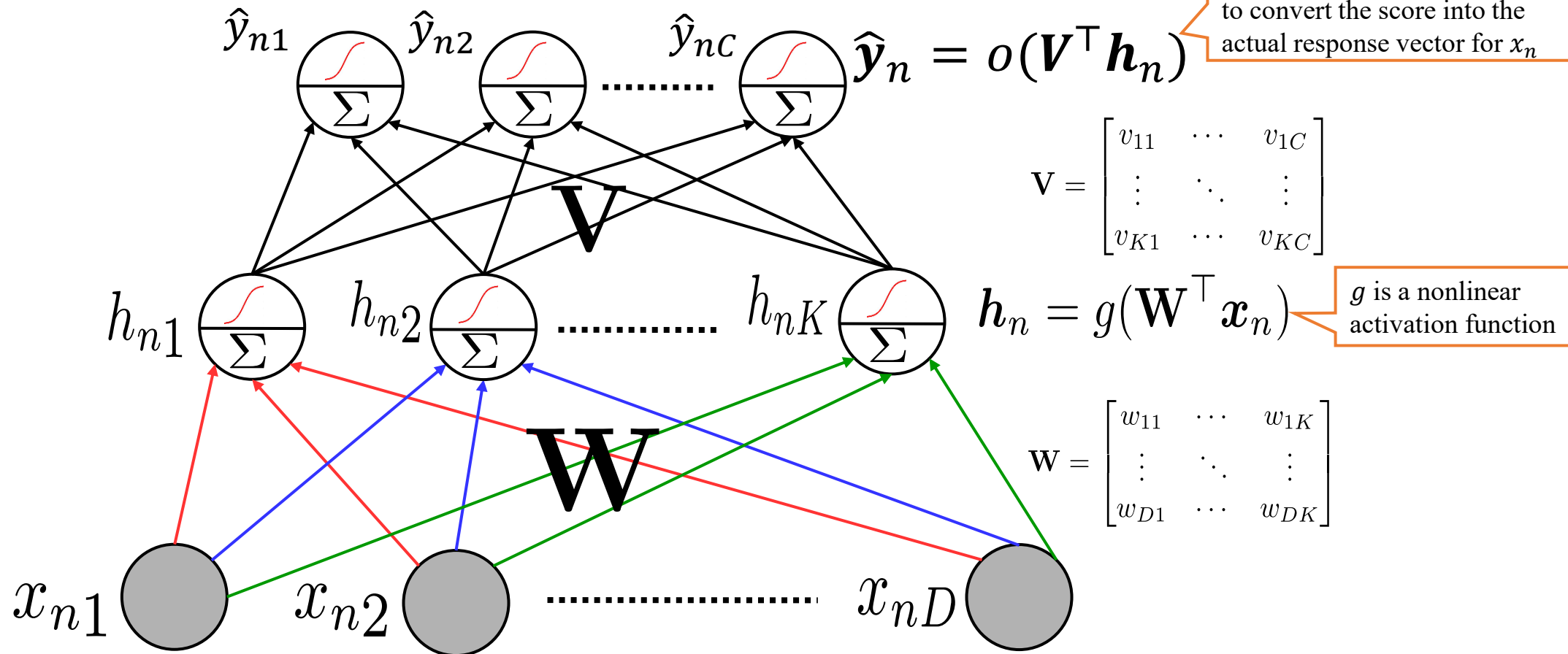
Single Hidden Layer and Single Output

- One hidden layer with K nodes and a single output (e.g., scalar-valued regression or binary classification)



Single Hidden Layer and Multiple Outputs

- One hidden layer with K nodes and a vector of C outputs (e.g., vector-valued regression or multi-class classification or multi-label classification)

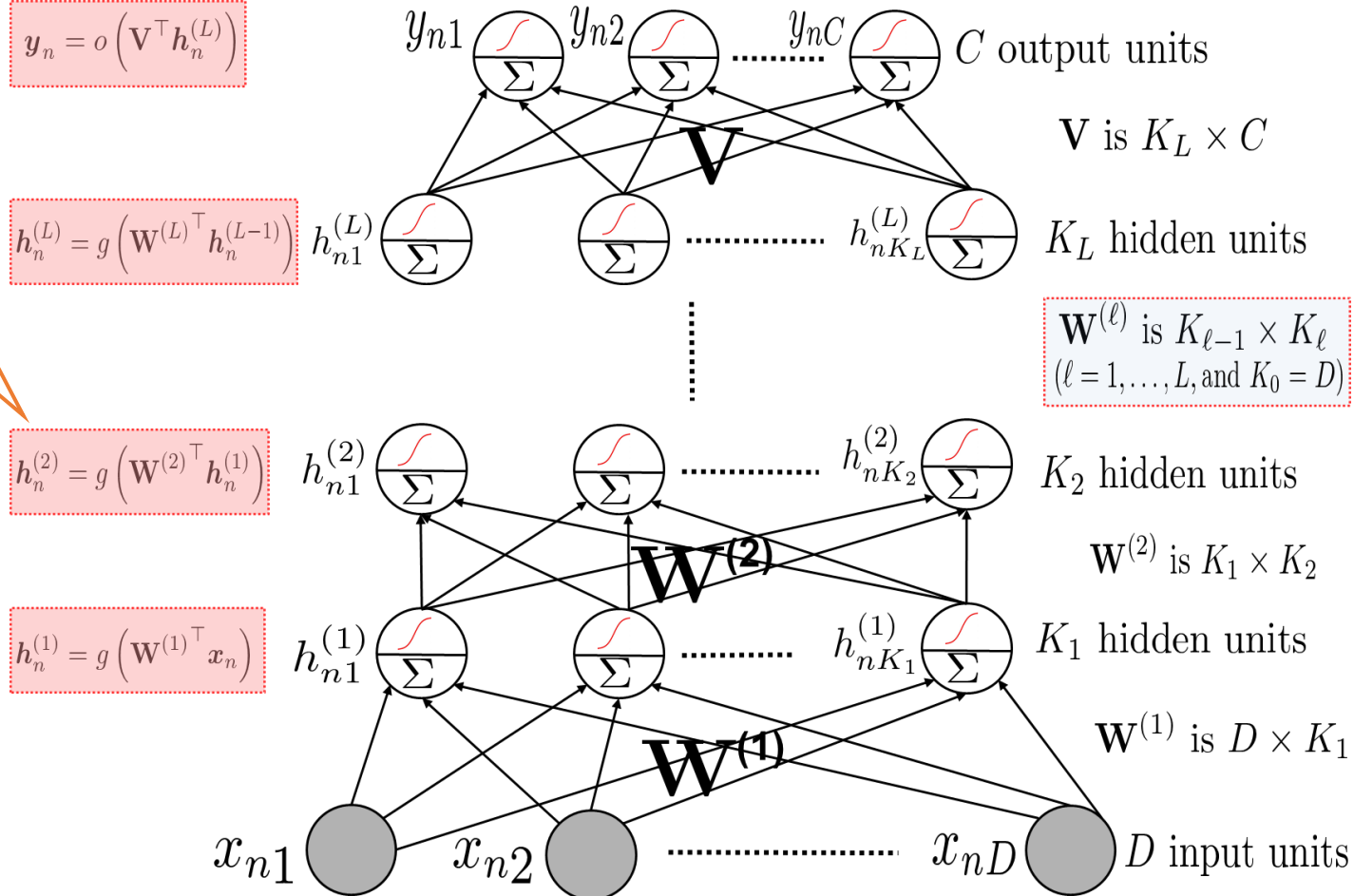


Multiple Hidden Layers (One/Multiple Outputs)

- Most general case: Multiple hidden layers with (with same or different number of hidden nodes in each) and a scalar or vector-valued output

Each hidden layer uses a nonlinear activation function g (essential, otherwise the network can't learn nonlinear functions and reduces to a linear model)

Note: Nonlinearity g is applied **element-wise** on its inputs so $h_n^{(\ell)}$ has the same size as vector $W^{(\ell)} h_n^{(\ell-1)}$



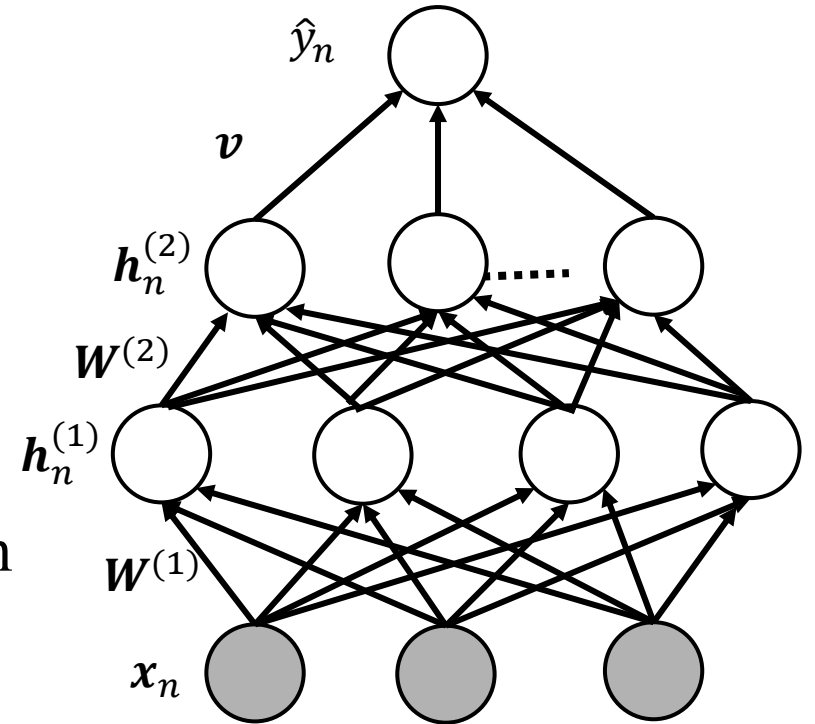
The Bias Term

- Each layer's pre-activations $\mathbf{z}_n^{(\ell)}$ have an add bias term $\mathbf{b}^{(\ell)}$ (has the same size as $\mathbf{z}_n^{(\ell)}$ and $\mathbf{h}_n^{(\ell)}$) as well

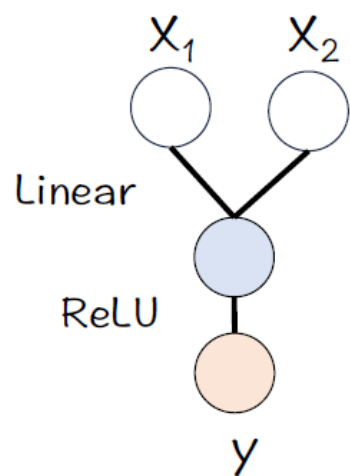
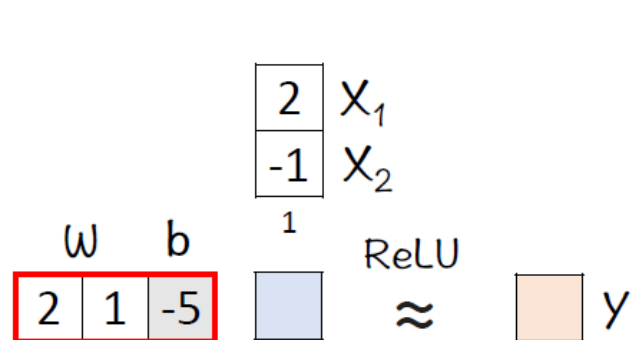
$$\mathbf{z}_n^{(\ell)} = \mathbf{W}^{(\ell)\top} \mathbf{x}_n^{(\ell-1)} + \mathbf{b}^{(\ell)}$$

$$\mathbf{h}_n^{(\ell)} = g(\mathbf{z}_n^{(\ell)})$$

- Bias term increases the expressiveness of the network and ensures that we have nonzero activations/pre-activations even if this layer's input is a vector of all zeros
- Note that the bias term is the same for all inputs (does not depend on n)
- The bias term $\mathbf{b}^{(\ell)}$ is also learnable

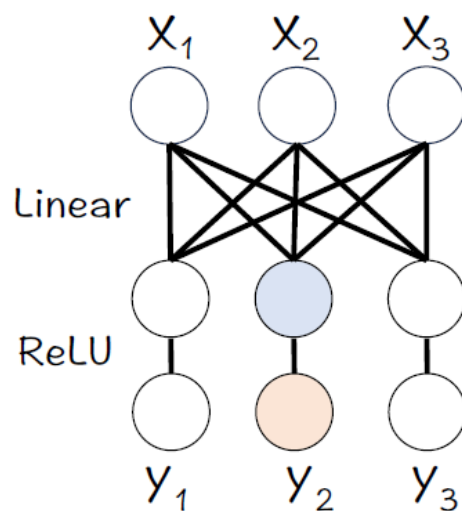
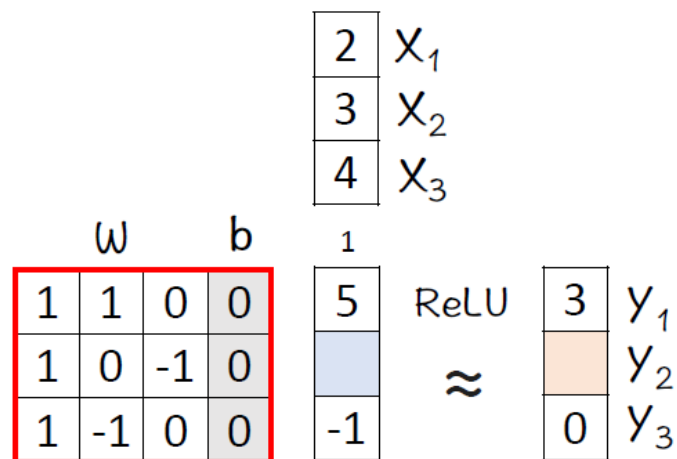


Exercises: Artificial Neuron



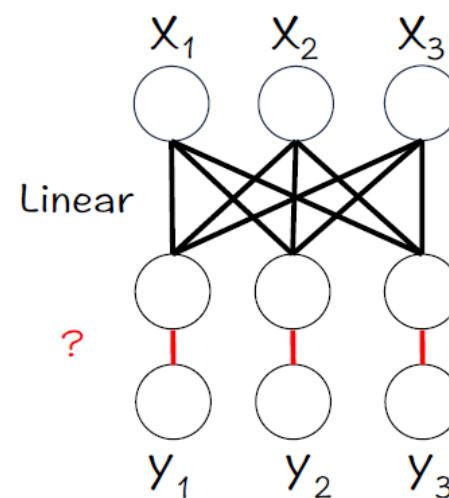
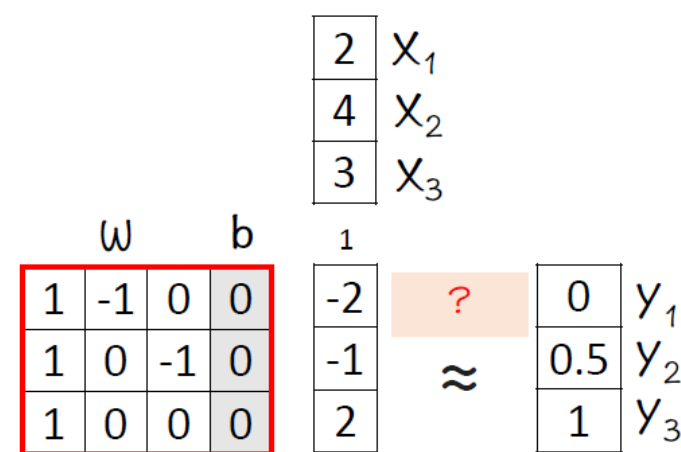
Exercise 1

-2
0



Exercise 2

-2
0

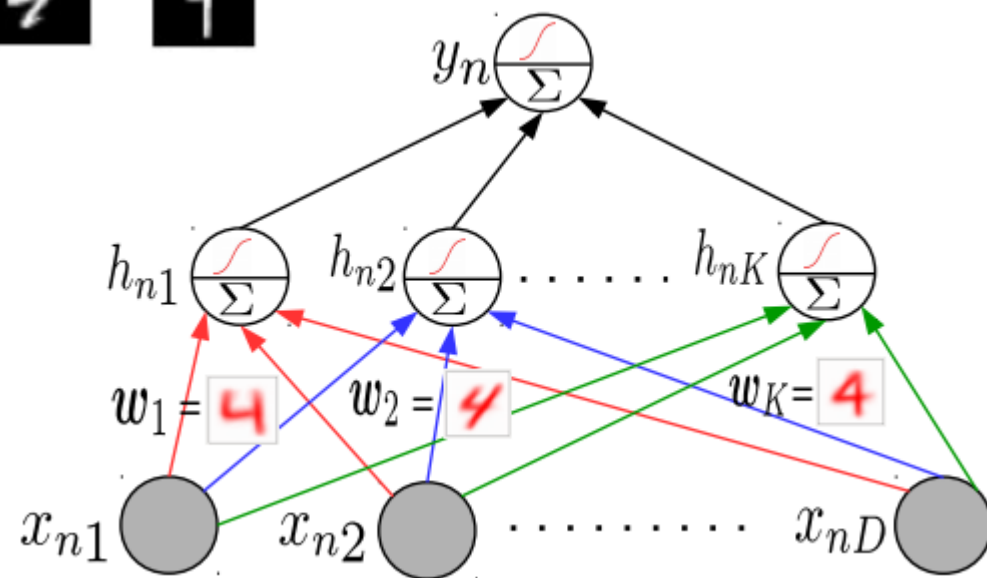
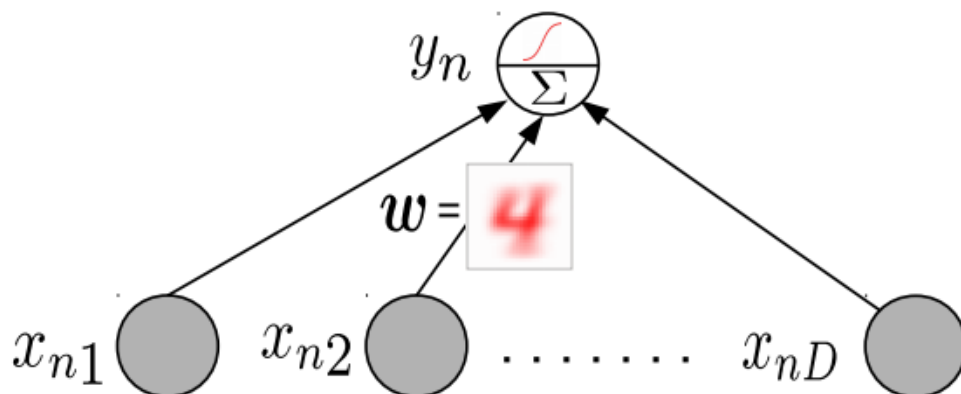


Exercise 3

Sigmoid

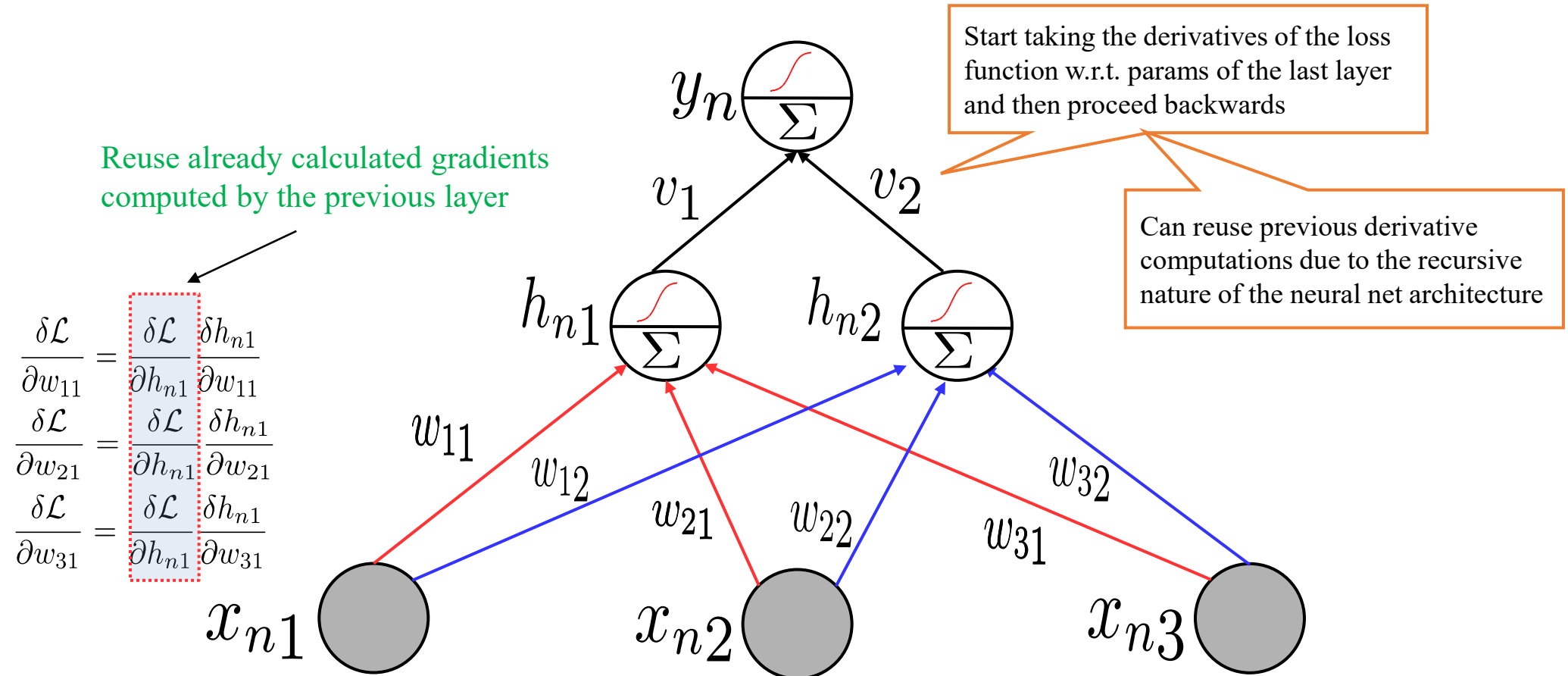
Why Neural Networks Work Better?

- Linear models tend to only learn the “average” pattern
 - E.g., Weight vector of a linear classification model represent average pattern of a class
- Deep models can learn multiple patterns (each hidden node can learn one pattern)
 - Thus, deep models can learn to capture more subtle variations that a simpler linear model



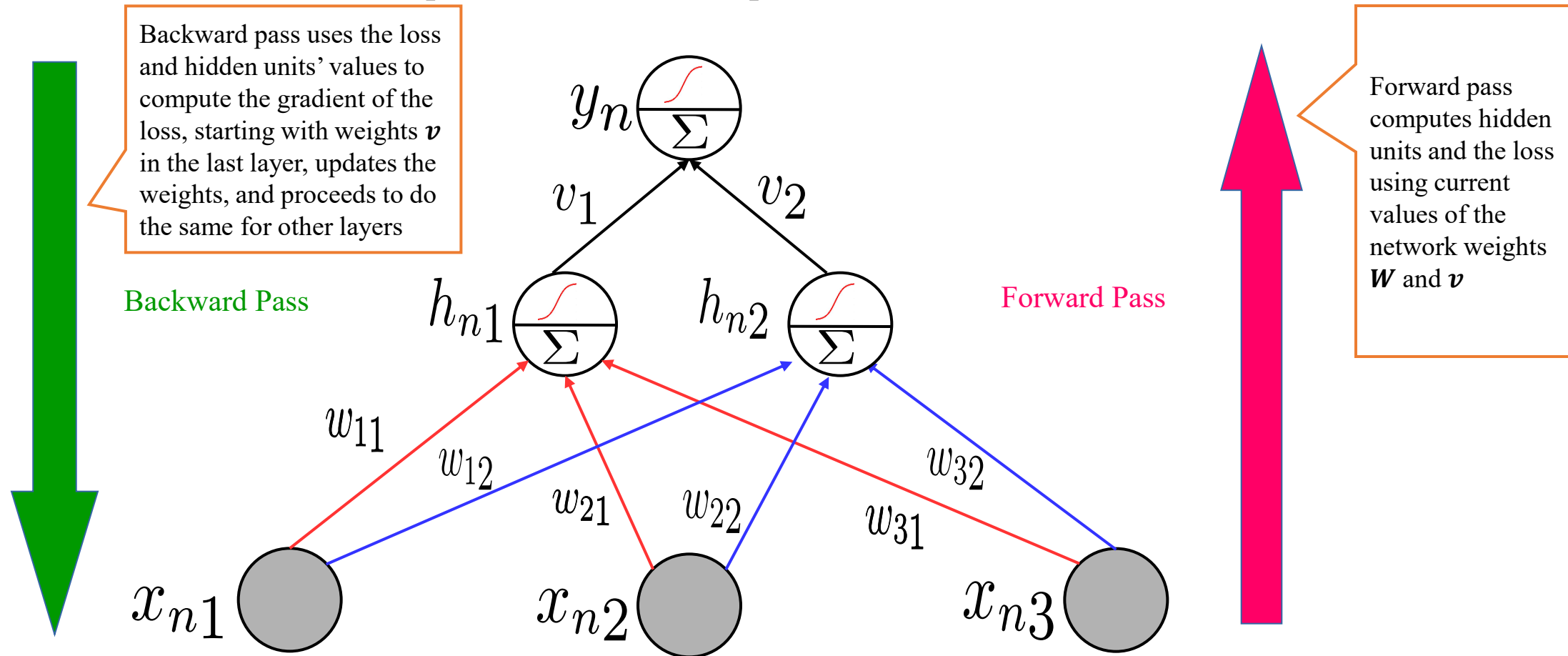
Backpropagation

- Backpropagation = Gradient descent using chain rule of derivatives
- Chain rule of derivatives: Example, if $y = f_1(x)$ and $x = f_2(z)$ then $\frac{\partial y}{\partial z} = \frac{\partial y}{\partial x} \frac{\partial x}{\partial z}$



Backpropagation

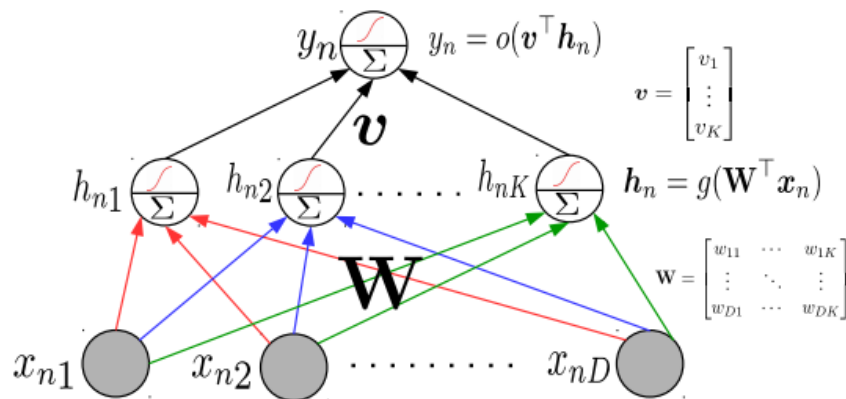
- Backprop iterates between a forward pass and a backward pass



- Software frameworks such as Tensorflow and PyTorch support this

Backpropagation through an example

Consider a single hidden layer MLP



Assuming regression ($o = \text{identity}$),
the loss function for this model

$$\begin{aligned}\mathcal{L} &= \frac{1}{2} \sum_{n=1}^N \left(y_n - \mathbf{v}^T \mathbf{h}_n \right)^2 \\ &= \frac{1}{2} \sum_{n=1}^N \left(y_n - \sum_{k=1}^K v_k h_{nk} \right)^2 \\ &= \frac{1}{2} \sum_{n=1}^N \left(y_n - \sum_{k=1}^K v_k g(\mathbf{w}_k^T \mathbf{x}_n) \right)^2\end{aligned}$$

- To use gradient methods for \mathbf{W} , \mathbf{v} , we need gradients.

- Gradient of \mathcal{L} w.r.t. \mathbf{v} is straightforward

$$\frac{\partial \mathcal{L}}{\partial v_k} = - \sum_{n=1}^N \left(y_n - \sum_{k=1}^K v_k g(\mathbf{w}_k^T \mathbf{x}_n) \right) h_{nk} = \sum_{n=1}^N \mathbf{e}_n h_{nk}$$

- Gradient of \mathcal{L} w.r.t. \mathbf{W} requires chain rule

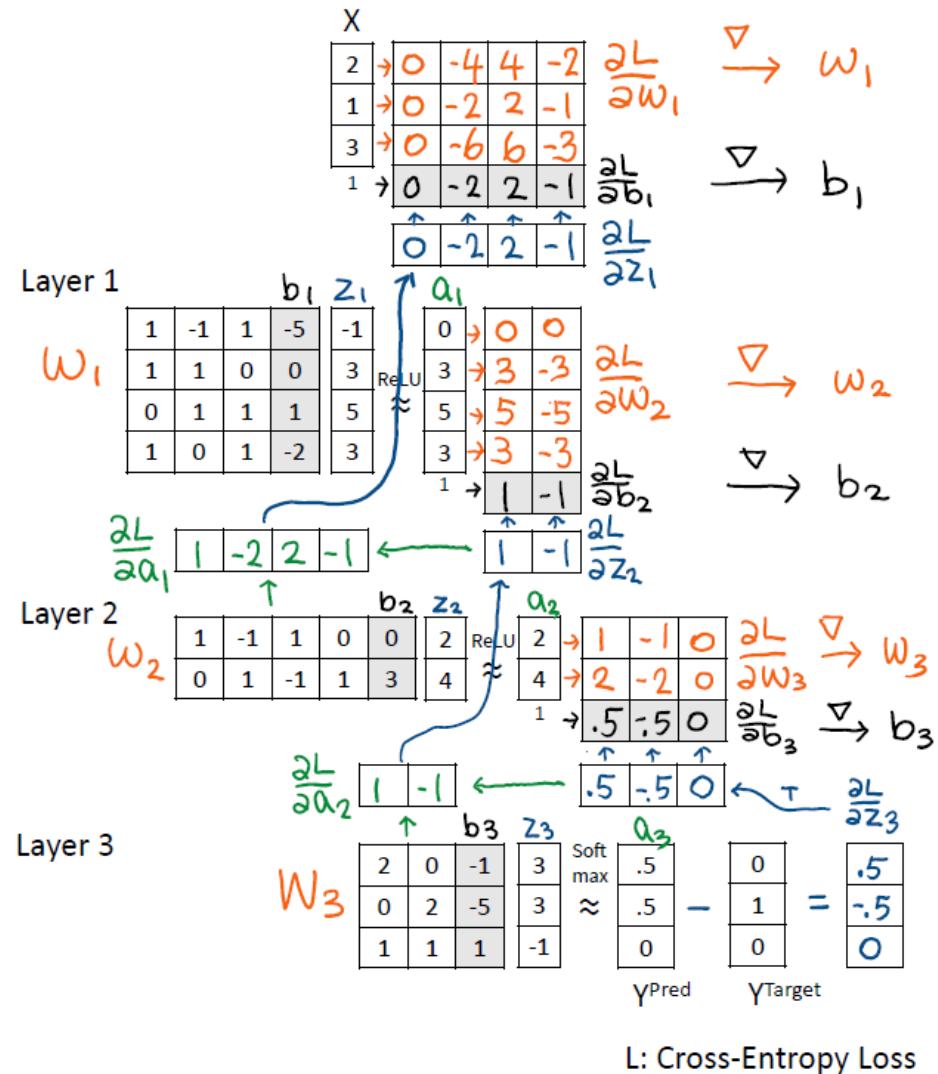
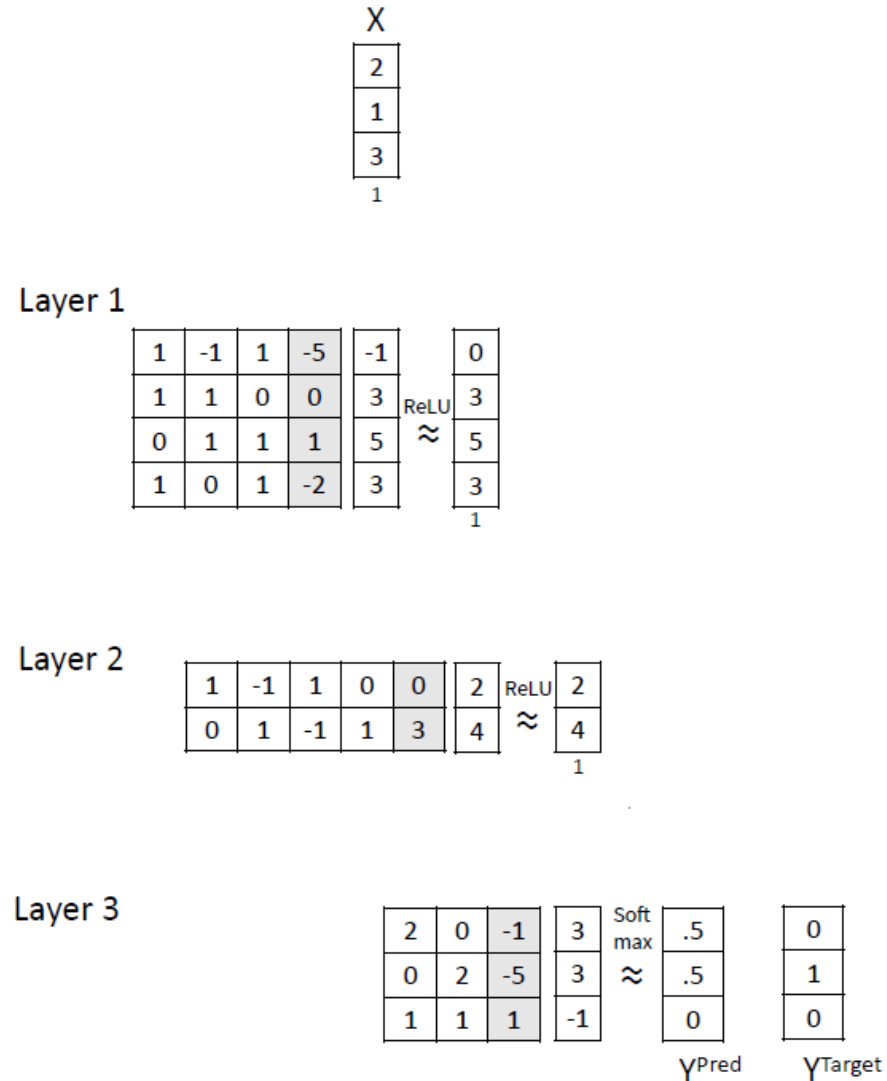
$$\frac{\partial \mathcal{L}}{\partial w_{dk}} = \sum_{n=1}^N \frac{\partial \mathcal{L}}{\partial h_{nk}} \frac{\partial h_{nk}}{\partial w_{dk}}$$

$$\frac{\partial \mathcal{L}}{\partial h_{nk}} = -(y_n - \sum_{k=1}^K v_k g(\mathbf{w}_k^T \mathbf{x}_n)) v_k = -\mathbf{e}_n v_k$$

$$\frac{\partial h_{nk}}{\partial w_{dk}} = g'(\mathbf{w}_k^T \mathbf{x}_n) x_{nd} \quad (\text{note: } h_{nk} = g(\mathbf{w}_k^T \mathbf{x}_n))$$

- **Forward prop** computes errors \mathbf{e}_n using current \mathbf{W} , \mathbf{v} .
Backprop updates NN params \mathbf{W} , \mathbf{v} using grad methods
- Backprop caches many of the calculations for reuse

Exercise: Backpropagation



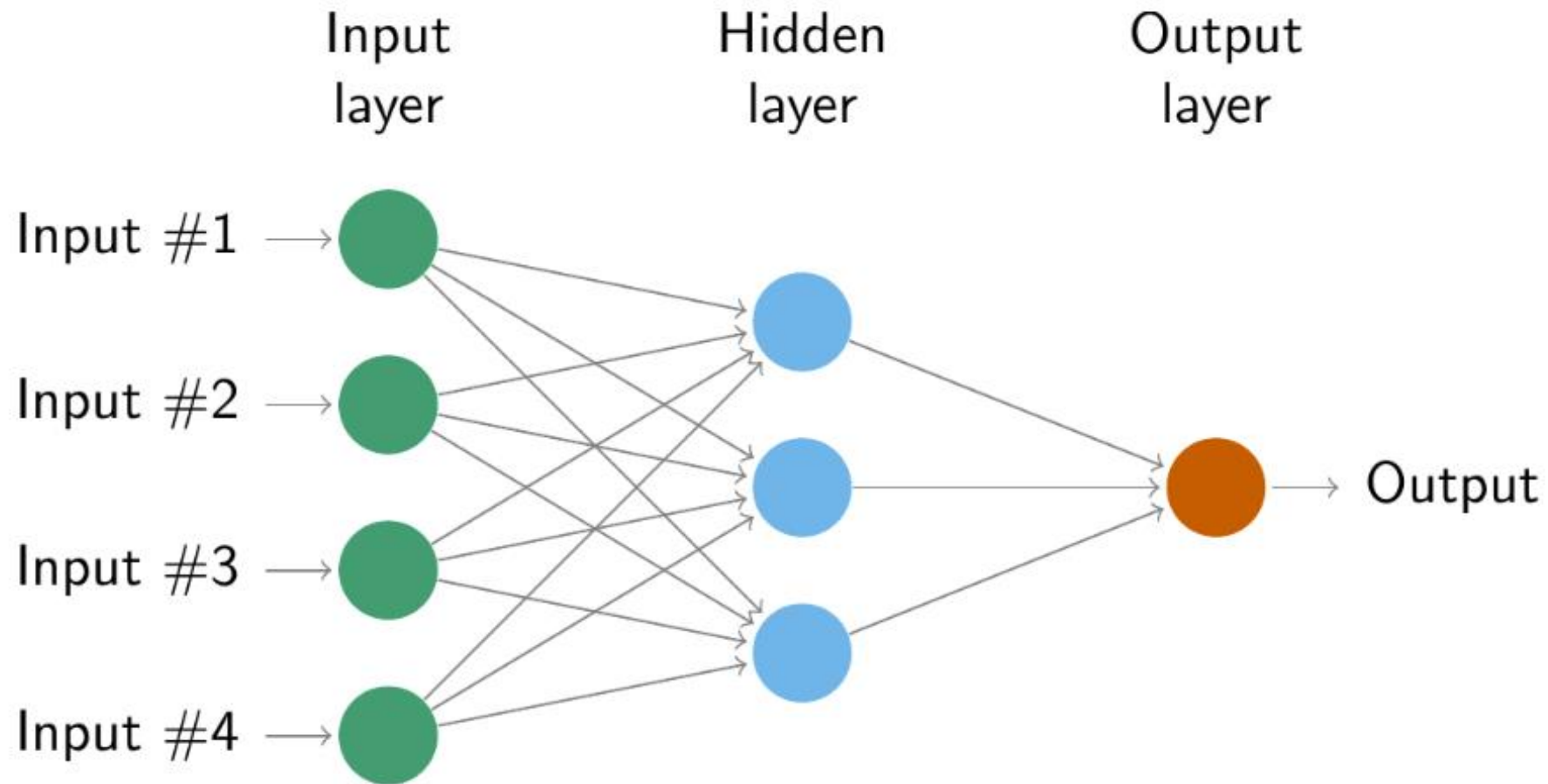
Exercise: Backpropagation

Steps:

1. **Forward Pass:** Given a multi layer perceptron (3 levels), an input vector X , predictions $Y^{Pred} = [0.5, 0.5, 0]$, and ground truth label $Y^{Target} = [0, 1, 0]$.
2. **Backpropagation:** Insert cells to hold our calculations.
3. **Layer 3** - Softmax (blue): Calculate $\partial L / \partial z_3$ directly using the simple equation: $Y^{Pred} - Y^{Target} = [0.5, -0.5, 0]$. This simple equation is the benefit of using Softmax and Cross Entropy Loss together.
4. **Layer 3** - Weights (orange) & Biases (black): Calculate $\partial L / \partial W_3$ and $\partial L / \partial b_3$ by multiplying $\partial L / \partial z_3$ and $[a_2 | 1]$.
5. **Layer 2** - Activations (green): Calculate $\partial L / \partial a_2$ by multiplying $\partial L / \partial z_3$ and W_3 .
6. **Layer 2** - ReLU (blue): Calculate $\partial L / \partial z_2$ by multiplying $\partial L / \partial a_2$ with 1 for positive values and 0 otherwise.
7. **Layer 2** - Weights (orange) & Biases (black): Calculate $\partial L / \partial W_2$ and $\partial L / \partial b_2$ by multiplying $\partial L / \partial z_2$ and $[a_1 | 1]$.
8. **Layer 1** - Activations (green): Calculate $\partial L / \partial a_1$ by multiplying $\partial L / \partial z_2$ and W_2 .
9. **Layer 1** - ReLU (blue): Calculate $\partial L / \partial z_1$ by multiplying $\partial L / \partial a_1$ with 1 for positive values and 0 otherwise.
10. **Layer 1** - Weights (orange) & Biases (black): Calculate $\partial L / \partial W_1$ and $\partial L / \partial b_1$ by multiplying $\partial L / \partial z_1$ and $[x | 1]$.
11. **Gradient Descent:** Update weights and biases (typically a learning rate is applied here).

Autoregressive Neural Network

Feed-forward Neural Network



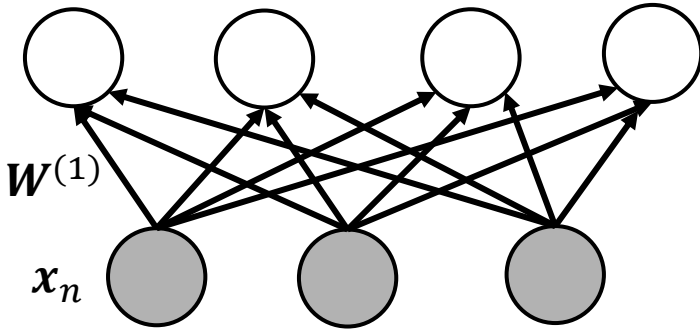
ARNN for Forecasting

- With time series data, lagged values of the time series can be used as inputs to a neural network, typically named as **Autoregressive neural network (ARNN)**
- $\text{ARNN}(p, k)$ is a feed-forward network with one hidden layer, with p lagged inputs and k nodes in the hidden layer.
- $\text{ARNN}(p, 0)$ model is equivalent to an $\text{ARIMA}(p, 0, 0)$ model, but without the restrictions on the parameters to ensure stationarity.
- With seasonal data, it is useful to also add the last observed values from the same season as inputs.
- For time series, the default is the optimal number of lags (according to the AIC) for a linear $\text{AR}(p)$ model. If k is not specified, it is set to $k = (p + 1)/2$ (rounded to the nearest integer).
- When it comes to forecasting, the network is applied iteratively.

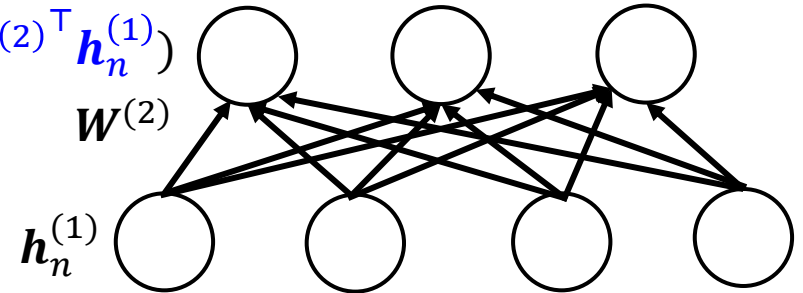
Limitations/Shortcomings of MLP & ARNN

- MLP uses fully connected layers defined by matrix multiplications + nonlinearity

$$h_n^{(1)} = g(W^{(1)\top} x_n)$$



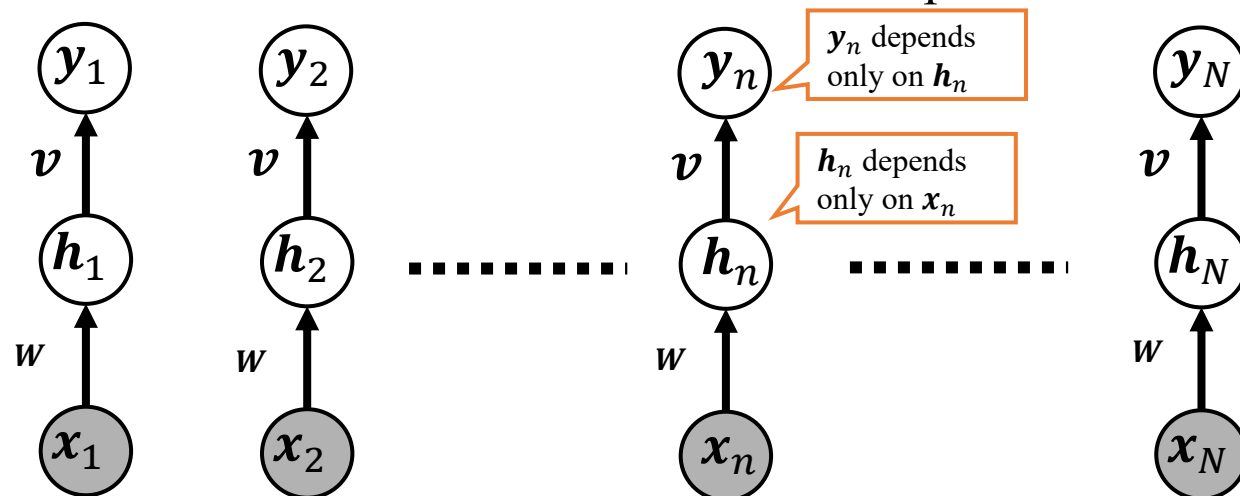
$$h_n^{(2)} = g(W^{(2)\top} h_n^{(1)})$$



- MLP **ignores structure** (e.g., spatial/sequential) in the inputs
 - Not ideal for data such as images, text, etc. which are flattened as vectors when used with MLP
- Fully connected nature of MLP requires massive number of weights
 - Recall that each layer is fully connected so each layer needs a massive number of weights!

Recurrent Connections in Deep Neural Networks

- Feedforward nets such as MLP assume independent observations

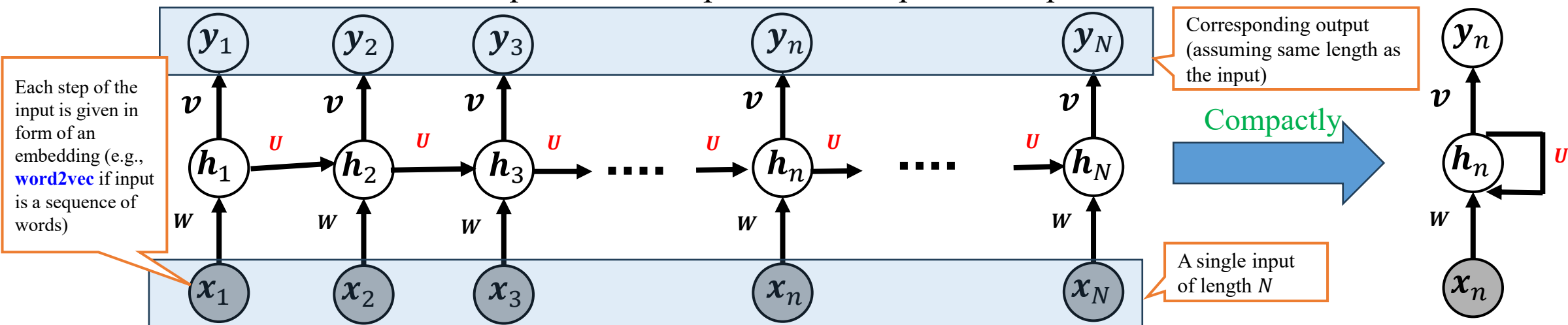


y_n depends only on h_n

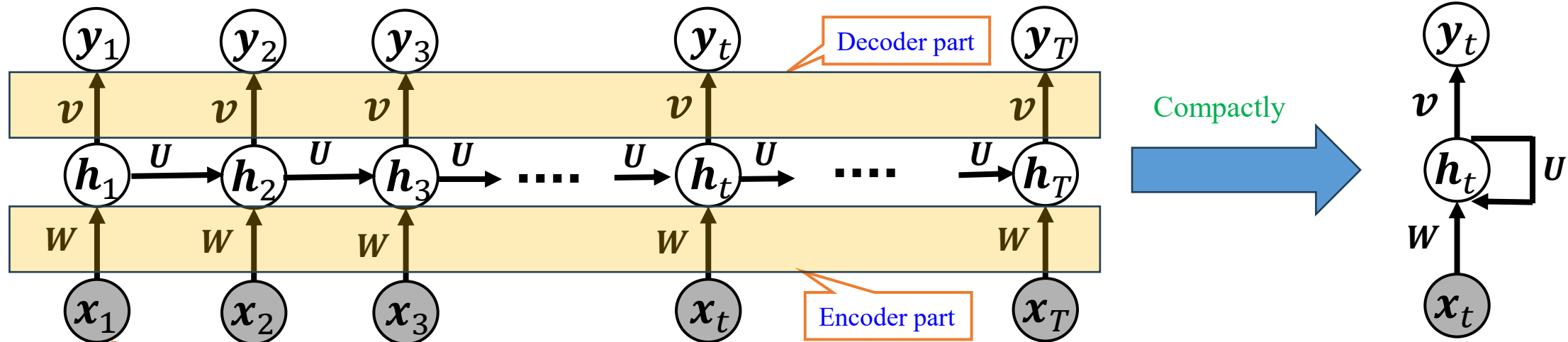
h_n depends only on x_n

Feedforward neural networks are not ideal when inputs $[x_1, x_2, \dots, x_N]$ and/or outputs $[y_1, y_2, \dots, y_N]$ represent sequential data (e.g., sequence of words, video (sequence of frames), etc).

- A **recurrent structure** can be helpful if each input and/or output is a sequence



- RNNs are used when each input or output or both are **sequences of tokens**

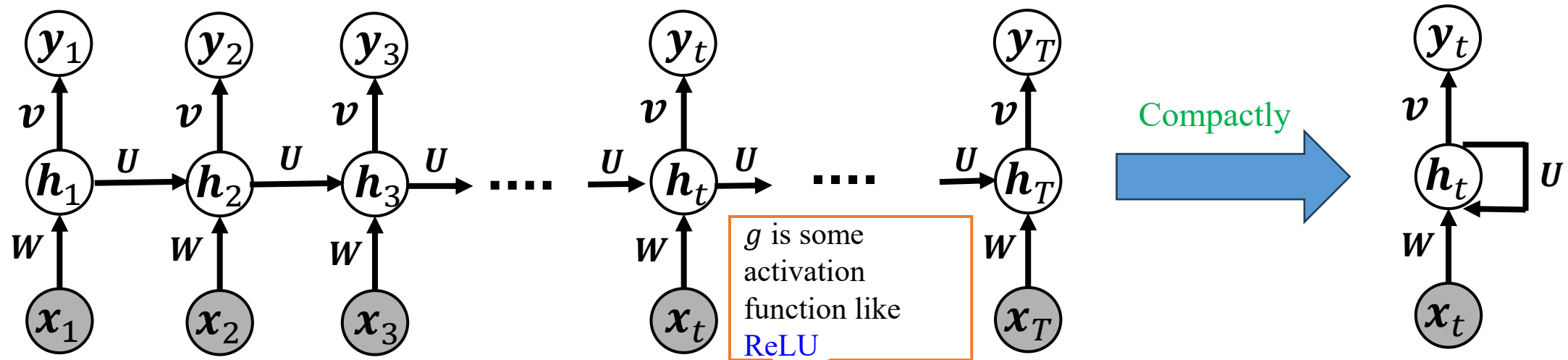


If the input is a word sequence, then each x_n represent the corresponding word's embedding (either a pre-computed word embedding like word2vec or a learned word embedding)

- Hidden state h_t is supposed to remember everything up to time $t - 1$. However, in practice, RNNs have difficulties remembering the distant past
 - Variants such as LSTM, GRU, etc mitigate this issue to some extent
- Slow processing is another major issue (e.g., can't compute h_t before computing h_{t-1})

Recurrent Neural Networks

- A basic RNN's architecture (assuming input and output sequence have same lengths)
- RNN has three sets of weights W, U, v



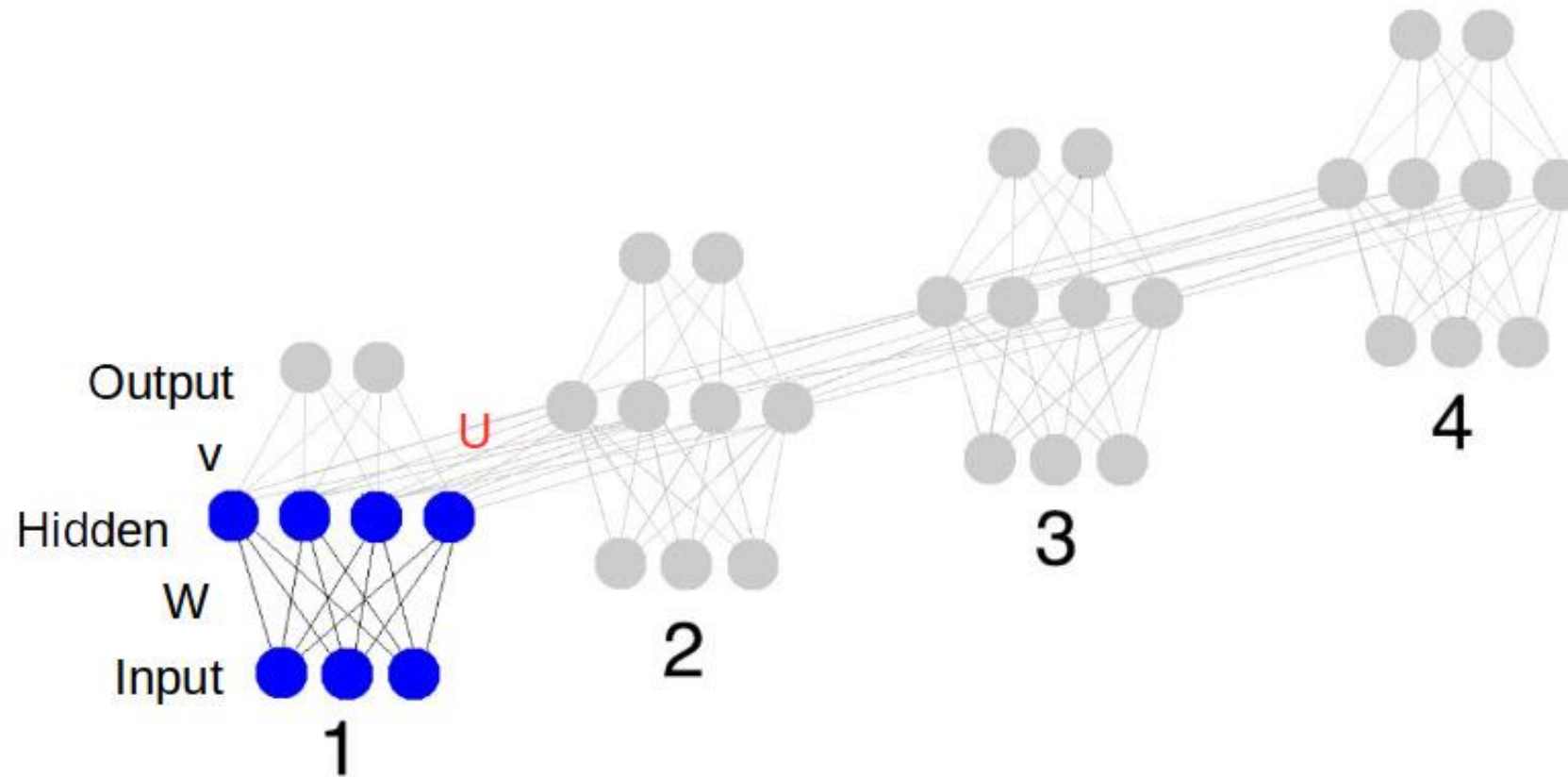
- W and U model how h_t at step t is computed: $h_t = g(Wx_t + Uh_{t-1})$
- v models the hidden layer to output mapping, e.g., $y_t = o(vh_t)$
- **Important:** Same W, U, v are used at all steps of the sequence (weight sharing)

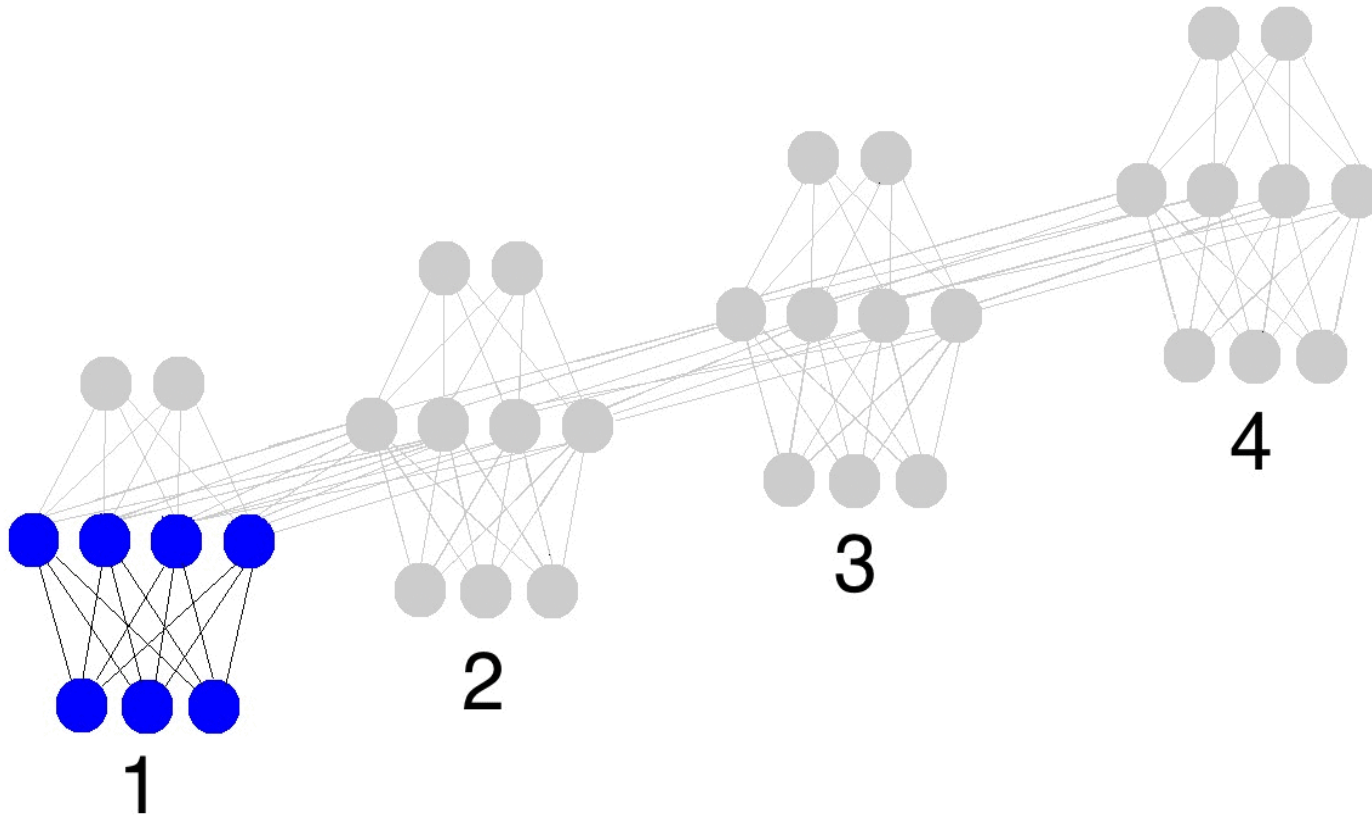
Given in form of an embedding
(e.g., word embedding if x_1 is a word)

o depends on the nature of y_t . If it is categorical then o can be **softmax**

Recurrent Neural Nets (RNN)

- A more “micro” view of RNN (the transition matrix U connects the hidden states across observations, propagating information along the sequence)





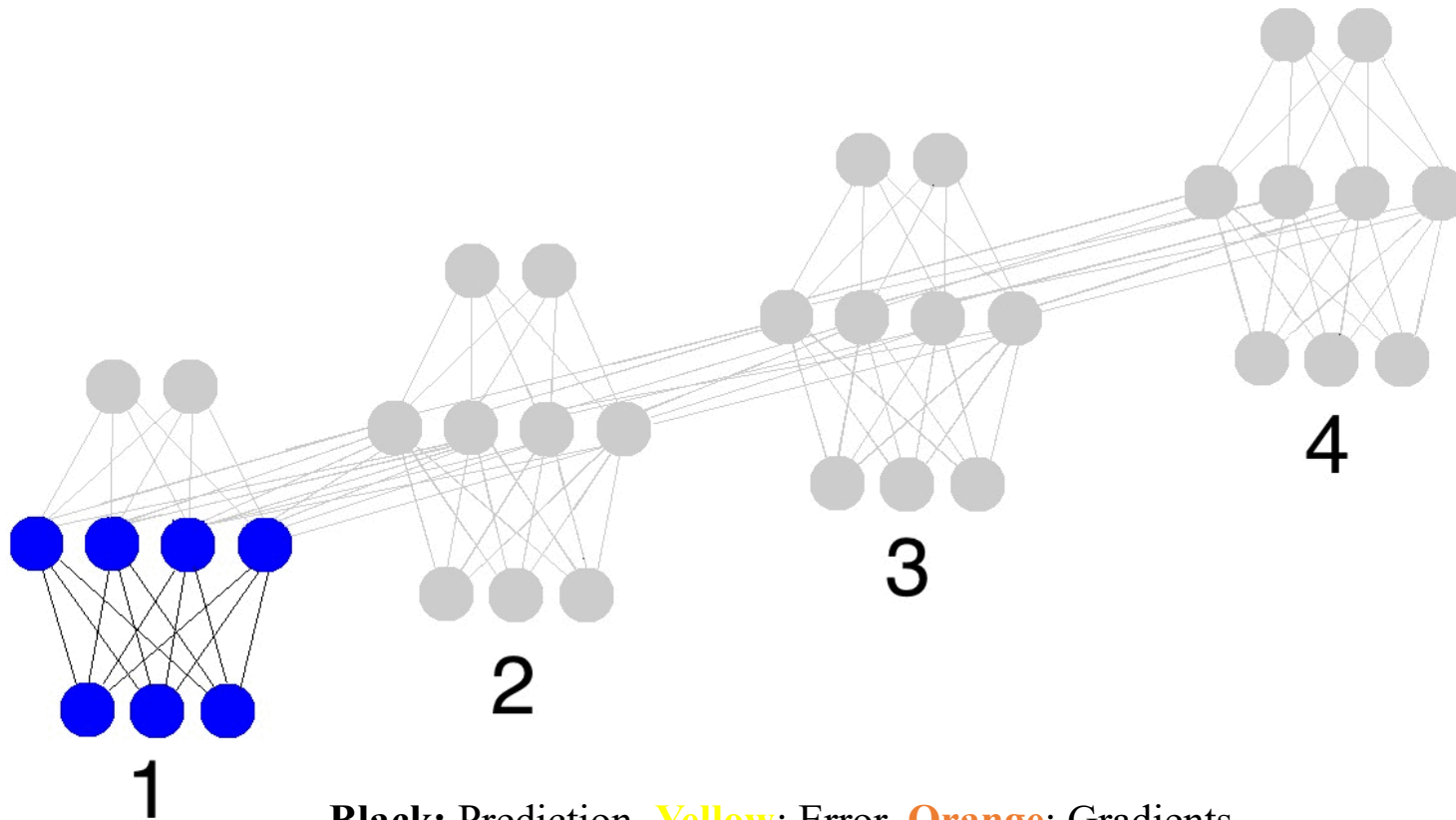
MakeAGIF.com

Workflow of RNN:

- The gif above reflects the magic of recurrent networks.
- It depicts 4 timesteps. The first is exclusively influenced by the input data.
- The second one is a mixture of the first and second inputs. This continues on.
- You should recognize that, in some way, network 4 is "full".
- Presumably, timestep 5 would have to choose which memories to keep and which ones to overwrite.
- This is very real. It's the notion of memory "capacity".
- As you might expect, bigger layers can hold more memories for a longer period of time.

Training RNN

- Trained using **Backpropagation Through Time** (forward propagate from step 1 to end, and then backward propagate from end to step)
- Think of the time-dimension as another hidden layer and then it is just like standard backpropagation for feedforward neural nets

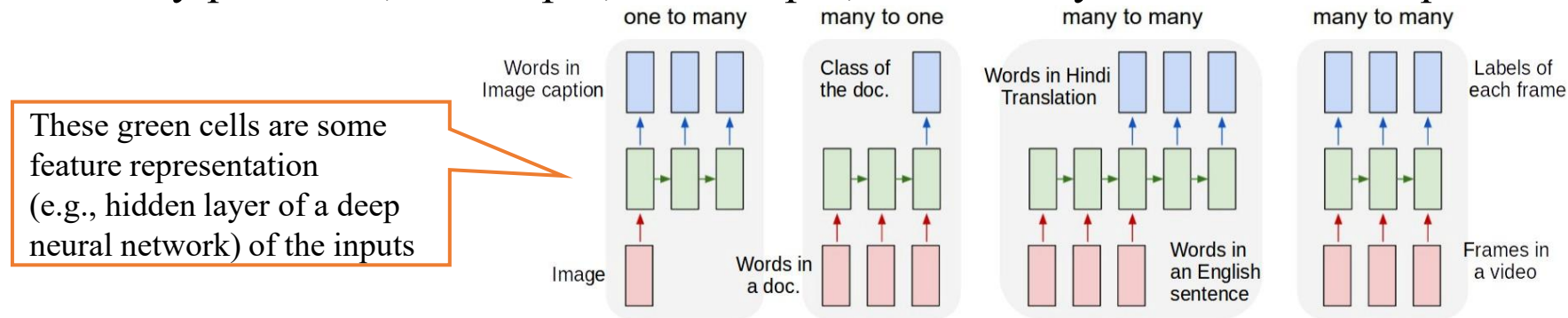


Black: Prediction, **Yellow:** Error, **Orange:** Gradients

- They learn by fully propagating forward from 1 to 4 (through an entire sequence of arbitrary length), and then backpropagating all the derivatives from 4 back to 1.
- You can also pretend that it's just a funny shaped normal neural network, except that we're re-using the same weights (synapses 0,1,and h) in their respective places.
- Other than that, it's normal backpropagation.

Pic source: <https://iamtrask.github.io/>

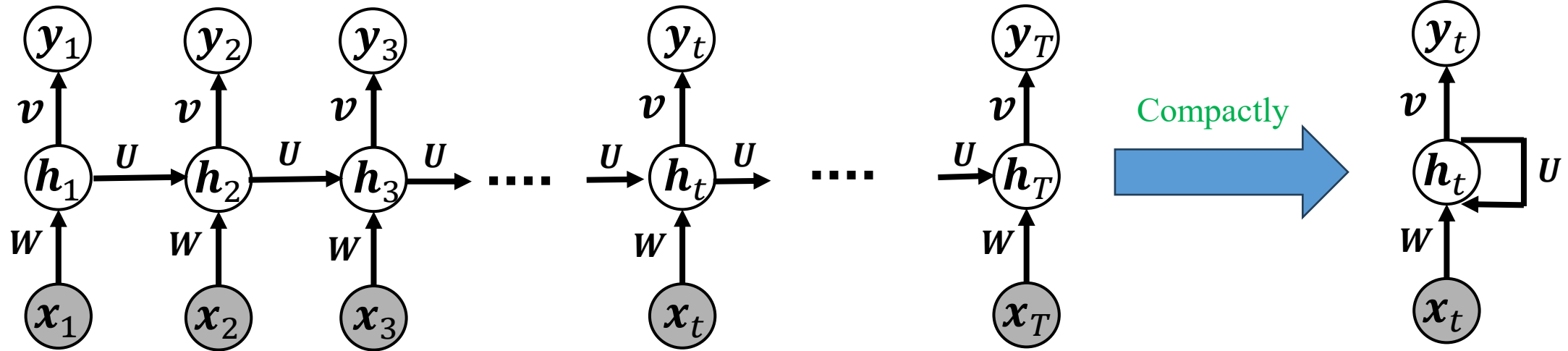
- In many problems, each input, each output, or both may be in form of sequences



- Different inputs or outputs need not have the same length
- Some examples of prediction tasks in such problems
 - **Image captioning:** Input is image (not a sequence), output is the caption (word sequence)
 - **Document classification:** Input is a word sequence, output is a categorical label
 - **Machine translation:** Input is a word sequence, output is a word sequence (in different language)
 - **Stock price prediction:** Input is a sequence of stock prices, output is its predicted price tomorrow
 - No input – just output (e.g., **generation** of random but plausible-looking text)

RNNs: Long Distant Past is Hard to Remember

- The hidden layer nodes h_t are supposed to summarize the past up to time $t - 1$

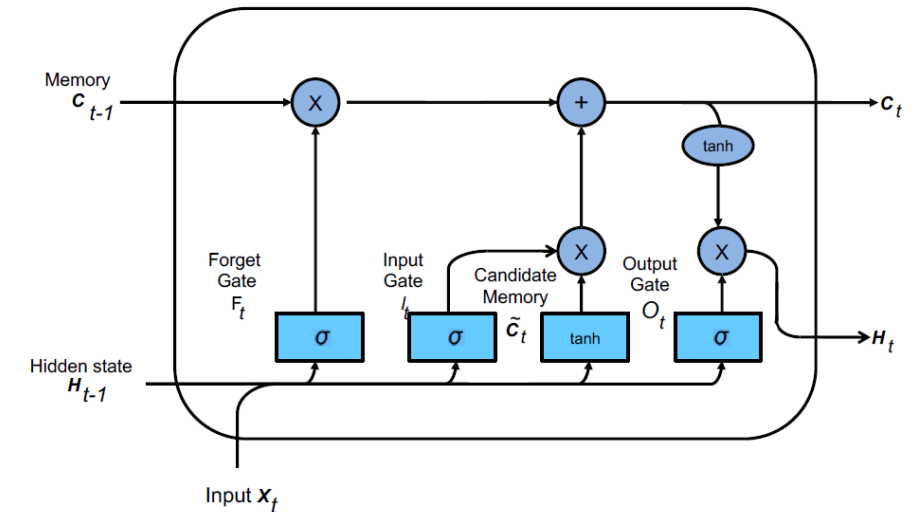
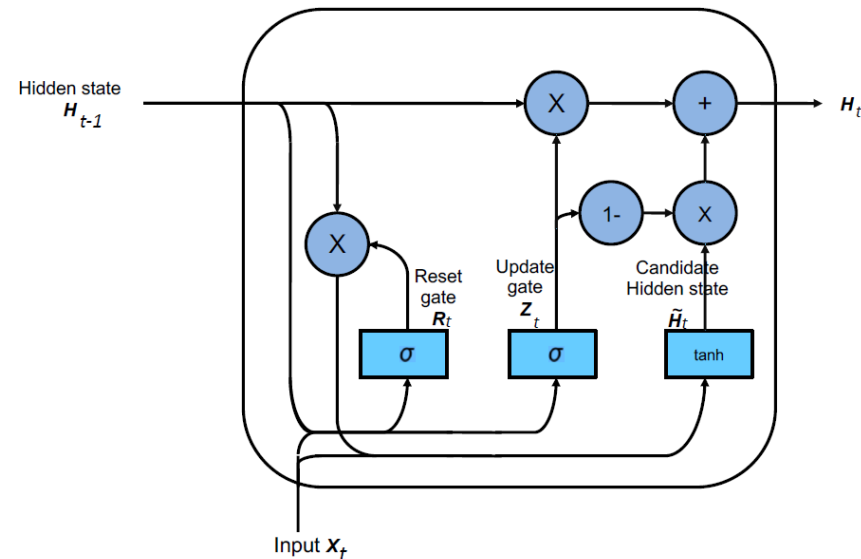


- In theory, they should. In practice, they can't. Some reasons
 - Vanishing gradients along the sequence too (due to repeated multiplications)
 - past knowledge gets “diluted”
 - Hidden nodes also have limited capacity because of their finite dimensionality
- Various extensions of RNNs have been proposed to address forgetting
 - Gated Recurrent Units (GRU), Long Short Term Memory (LSTM)

- Essentially an RNN, except that the hidden states are computed differently
- Recall that RNN computes the hidden states as

$$\mathbf{h}_t = \tanh(\mathbf{W}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1})$$

- For RNN: State update is multiplicative (weak memory and gradient issues)
- GRU and LSTM contain specialized units and “memory” which modulate what/how much information from the past to retain/forget



- In contrast, LSTM maintains a “context” \mathbf{C}_t and computes hidden states as

$$\hat{\mathbf{C}}_t = \tanh(\mathbf{W}^c \mathbf{x}_t + \mathbf{U}^c \mathbf{h}_{t-1}) \quad (\text{“local” context, only up to immediately preceding state})$$

$$i_t = \sigma(\mathbf{W}^i \mathbf{x}_t + \mathbf{U}^i \mathbf{h}_{t-1}) \quad (\text{how much to take in the local context})$$

$$f_t = \sigma(\mathbf{W}^f \mathbf{x}_t + \mathbf{U}^f \mathbf{h}_{t-1}) \quad (\text{how much to forget the previous context})$$

$$o_t = \sigma(\mathbf{W}^o \mathbf{x}_t + \mathbf{U}^o \mathbf{h}_{t-1}) \quad (\text{how much to output})$$

$$\mathbf{C}_t = \mathbf{C}_{t-1} \odot f_t + \hat{\mathbf{C}}_t \odot i_t \quad (\text{a modulated additive update for context})$$

$$\mathbf{h}_t = \tanh(\mathbf{C}_t) \odot o_t \quad (\text{transform context into state and selectively output})$$

- Note: \odot represents elementwise vector product. Also, state updates now additive, not multiplicative. Training using backpropagation through time.
- Many variants of LSTM exists, e.g., using \mathbf{C}_t in local computations, Gated Recurrent Units (GRU), etc. Mostly minor variations of basic LSTM above.

Exercise: RNN

Recurrent Neural Network (RNN)

Input Sequence X

3	4	5	6
---	---	---	---

Parameters A

1	-1
1	1

 B

1
2

 C

-1	1
----	---

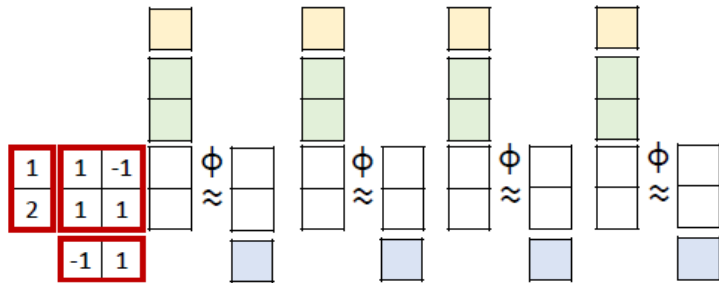
Activation Function ϕ : ReLU

Hidden States H_0

0
0

Output Sequence Y

--	--	--	--



Recurrent Neural Network (RNN)

Input Sequence X

3	4	5	6
---	---	---	---

Parameters A

1	-1
1	1

 B

1
2

 C

-1	1
----	---

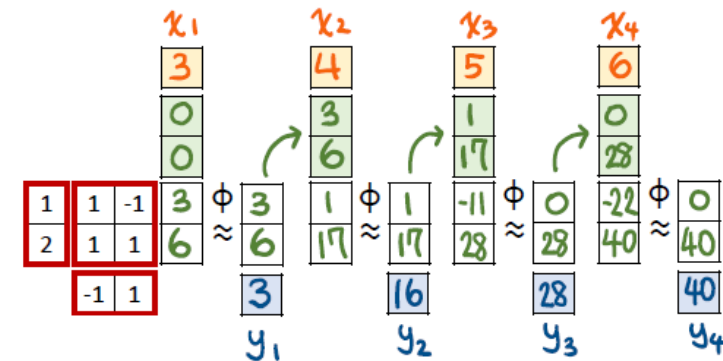
Activation Function ϕ : ReLU

Hidden States H_0

0
0

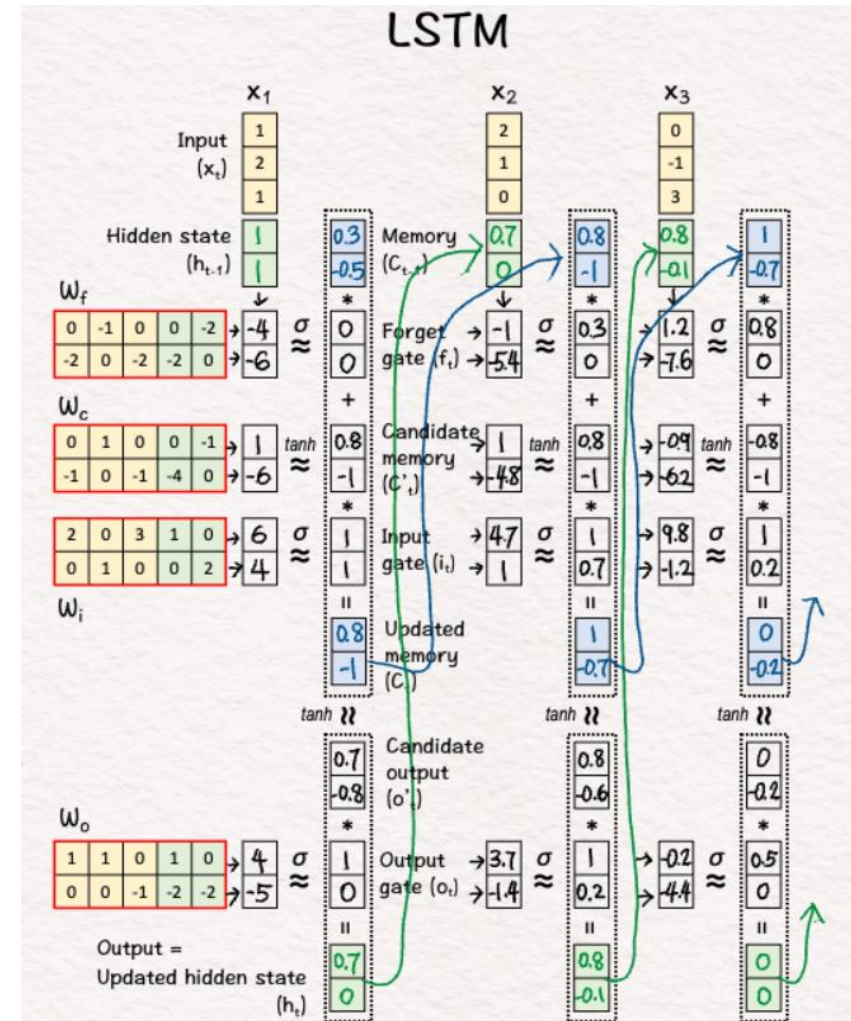
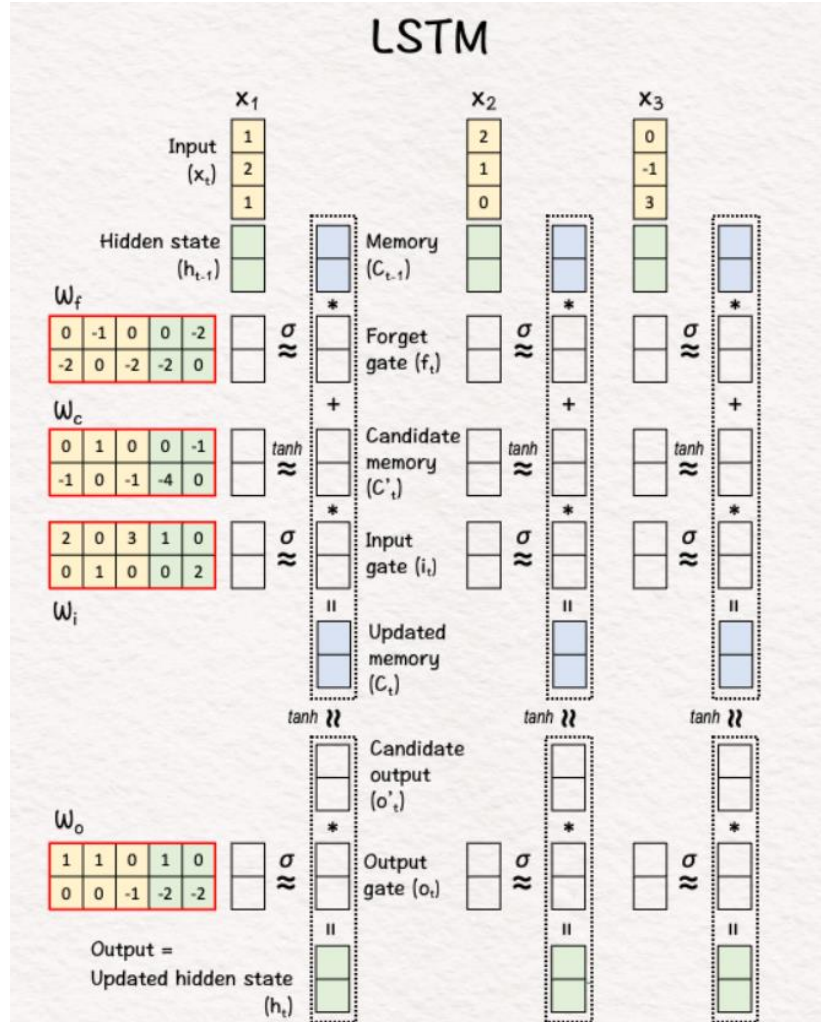
Output Sequence Y

--	--	--	--



Exercise: LSTM

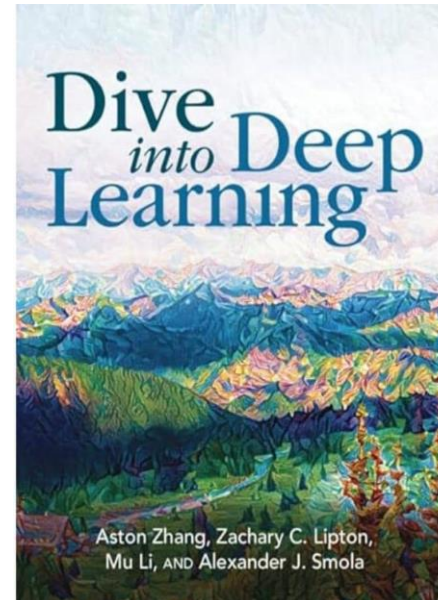
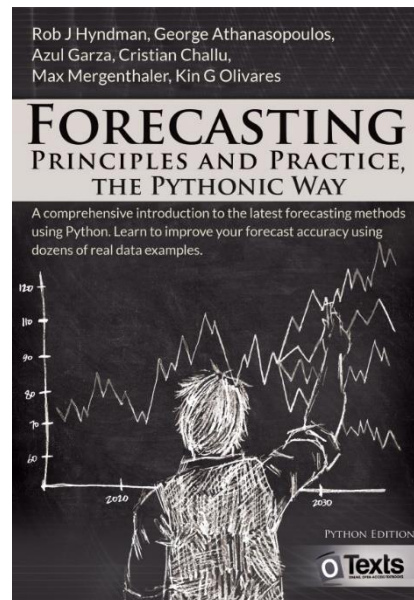
Initialize: Randomly set the previous hidden state h_0 to $[1, 1]$ and memory cells C_0 to $[0.3, -0.5]$



Pic source: <https://www.byhand.ai/>

HAPPY FORECASTING

“A good forecaster is not smarter than everyone else, he merely has his ignorance better organised.”



Read Online: <https://otexts.com/fpppy/> and <https://d2l.ai/>



If you're brave enough to
say goodbye, life will reward
you with a new hello.

Paulo Coelho

