

Natural Language Processing with Recurrent Neural Networks

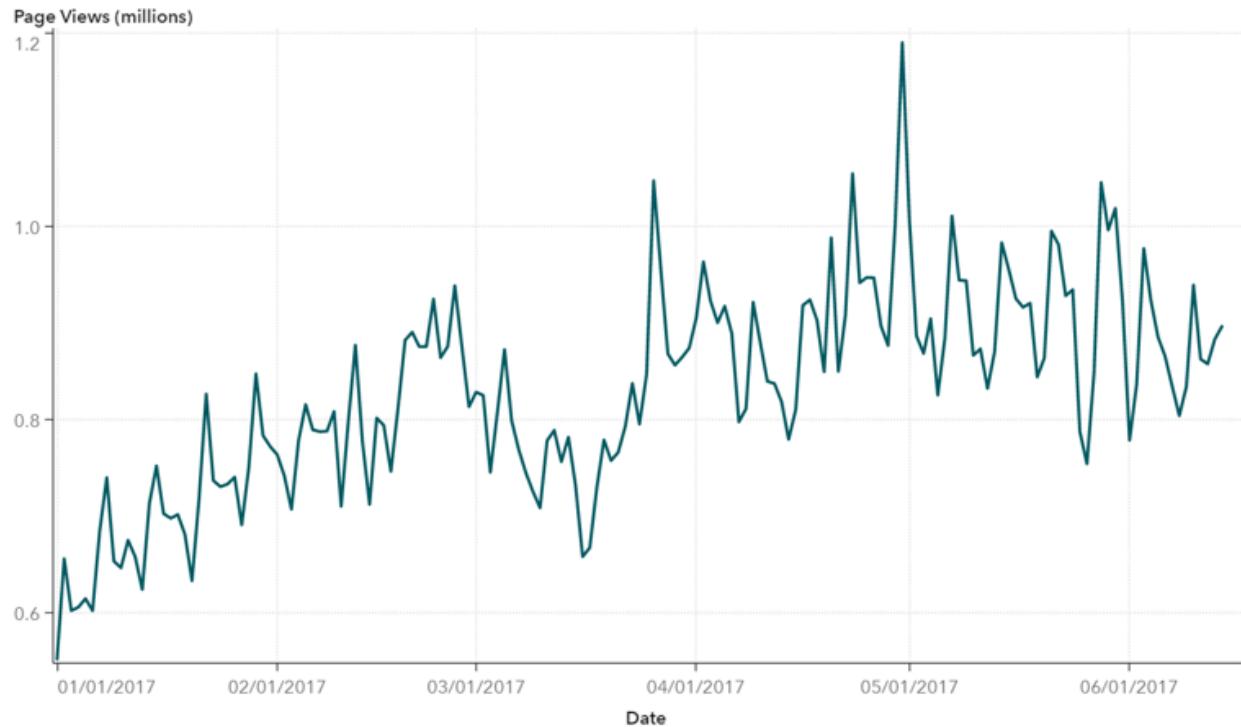
The structure of the lecture:

1. Why RNNs & NLP matter
2. RNNs: An introduction
3. RNNs Under the Hood (and architectural variations)
4. NLP: Why we use RNNs
5. Classifying sentiment: A coded example
6. Further Reading

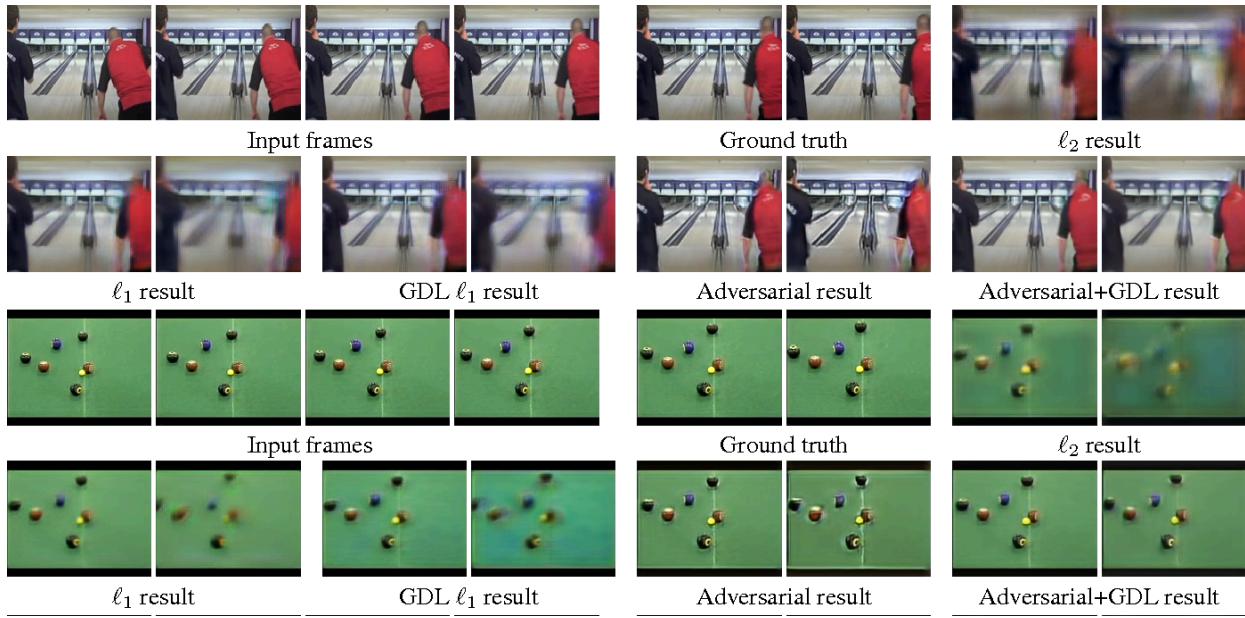
1 Why RNNs & NLP matter

Recurrent Neural Networks are Neural Networks specifically designed to deal with **sequences** as input data, i.e. **observations repeated throughout time**.

Example 1: Prediction of future stock market values and trends



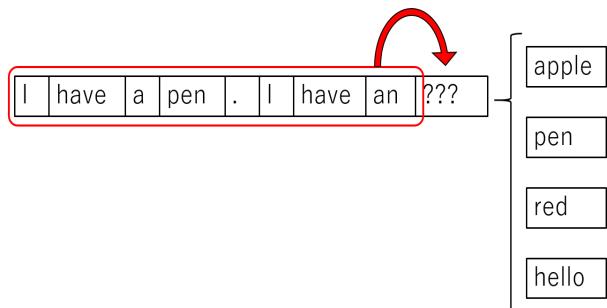
Example 2: Video prediction



🎥 Videos = sequences of images/frames

👉 Why not predicting the next image(s)?

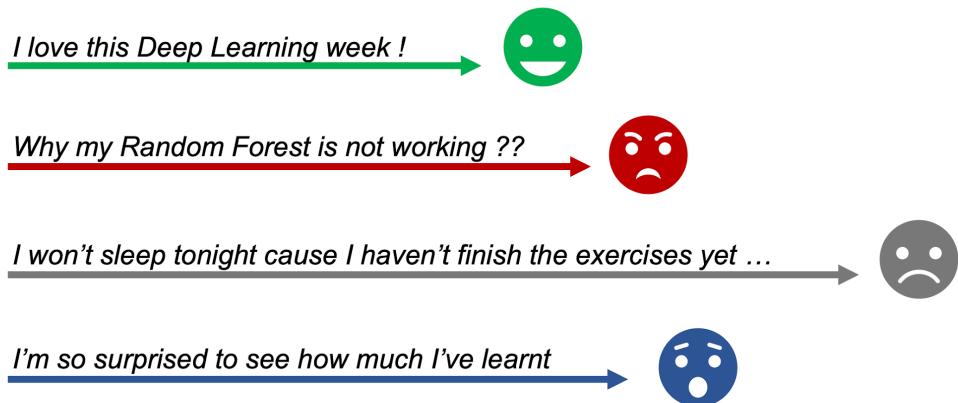
Example 3: Predicting the next word - Natural Language Processing



💡 Recurrent Networks are massively used for text!

NLP: Text classification, such as sentiment analysis

Classification depending on a word, a sentence, a paragraph, ...



The typical setting is **sentiment analysis**: Classify positive or negative sentences (but also happiness, sadness, joy, anger, ...).

Sequence to sequence models

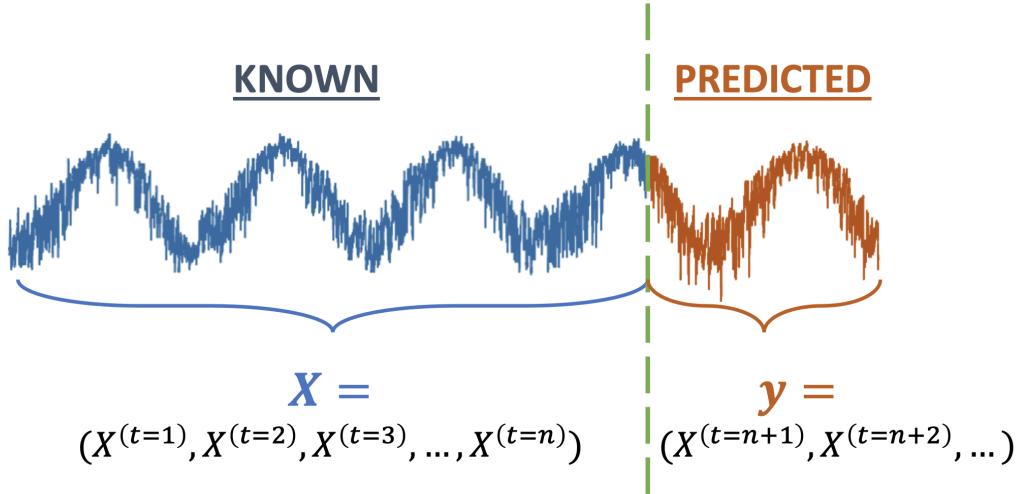
Given an input sequence, produce a corresponding output sequence. Typical application is language translation.

The screenshot shows the Google Translate interface. The input text is "Le wagon, a place where you will master deep learning". The output text is "Le wagon, un lieu où vous maîtriserez le deep learning". The interface includes language detection, document upload, and translation controls.

2 RNNs: An introduction

Inputs & Outputs:

During the Machine Learning lectures, you discovered Time Series as follows:



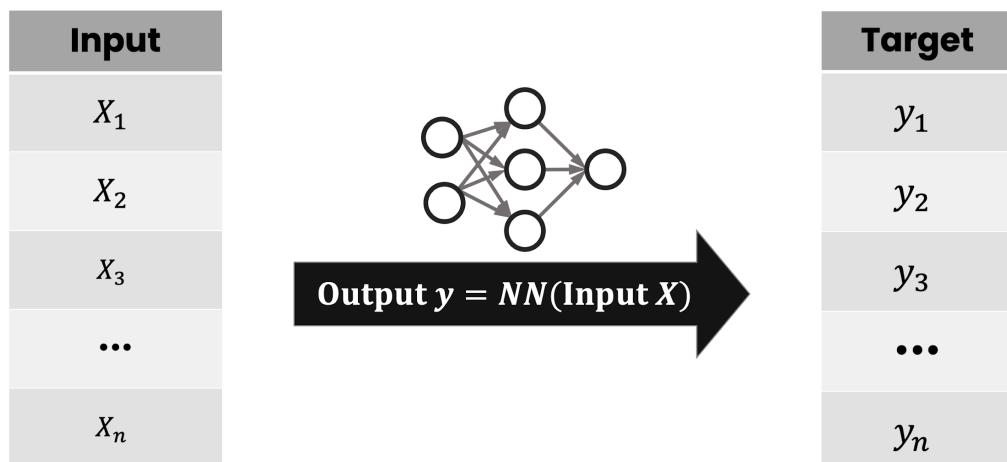
The algorithm would learn from the past values to predict the future values using temporal features such as the trend, the seasonality...

👉 Example: ARIMA models **recursively** predict the next data points **one after the other**

What follows is a slightly different approach...

General DL Framework

- **Training:** Many pairs of (
- X
- i
- ,
- y
- i
-) = (observation, target) are used to learn the patterns between the input and the output.



- **Prediction:** New input
- X

- n
- e
- w
- →
- New prediction/output
- y
- p
- r
- e
- d
-

! Each individual observation

X

i

can be of any type: scalar, vector, matrix (an image), ...

👉 X.shape = (n_ROWS, n_FEATURES)

What are our inputs for RNNs?

- **Training:** Each input
- X
- i
- consists of a **sequence of repeated observations** through time:

X

i

=

(

X

t

=

1

i

,

X

t

=

2

i

,

X

t

=

3

i

,

.

.

,

X

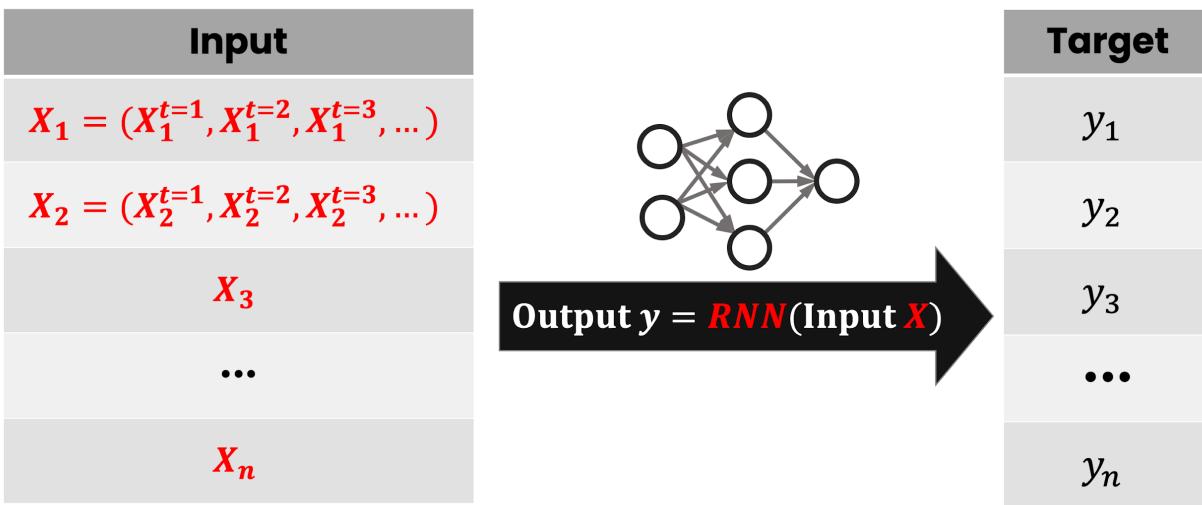
t

=

m

i

)



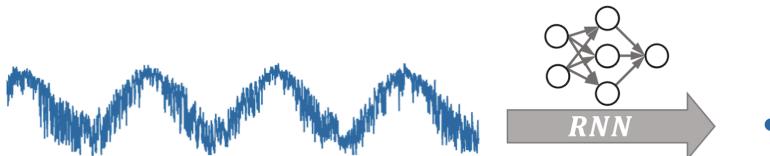
👉 `X.shape` = (n_SEQUENCES, n_OBSERVATIONS, n_FEATURES)

❗ The input tensor

`X`

has a new dimension, the **temporal dimension** which corresponds to the **number of observations through time**.

🔮 **Prediction: The next observation**

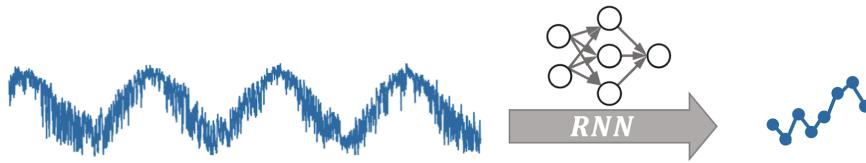


👉 Predict `rain` on day + 1

- `X.shape` = (None, 100 days, 1 feature)
- `y.shape` = (None, 1 next day, 1 feature)

ℹ️ None is a placeholder for the **number of sequences** given, it is equal to `batch_size`

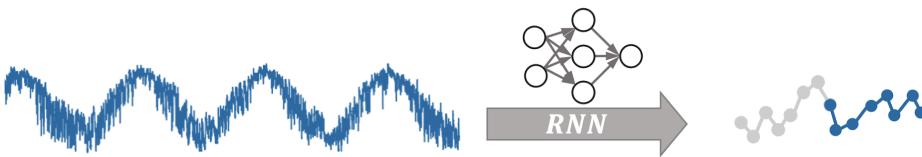
🔮 **Prediction: The next observation**



👉 Predict `rain` on [day + 1 ... day + 8]

- `X.shape = (None, 100 days, 1 feature)`
- `y.shape = (None, 8 next days, 1 feature)`

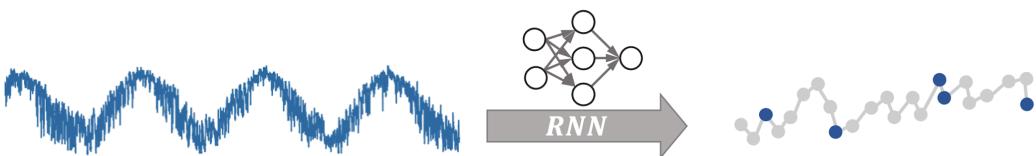
🔮 Prediction: The next observationS, but not necessarily right after the last seen value



👉 Predict `rain` on [day + 8 ... day + 16]

- `X.shape = (None, 100 days, 1 feature)`
- `y.shape = (None, 8 consecutive days in the future, 1 feature)`

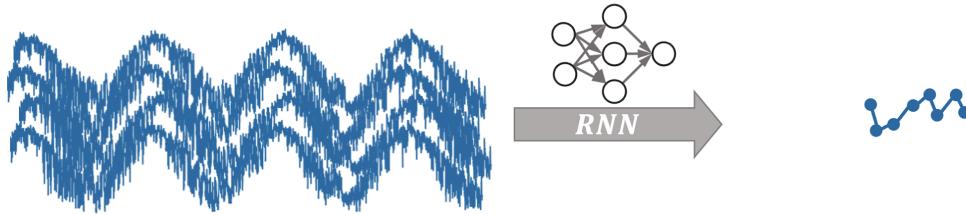
🔮 Prediction: The next observationS, but not necessarily consecutive



👉 Predict `rain` for 5 specific days in the future

- `X.shape = (None, 100 days, 1 feature)`
- `y.shape = (None, 5 specific days in the future, 1 feature)`

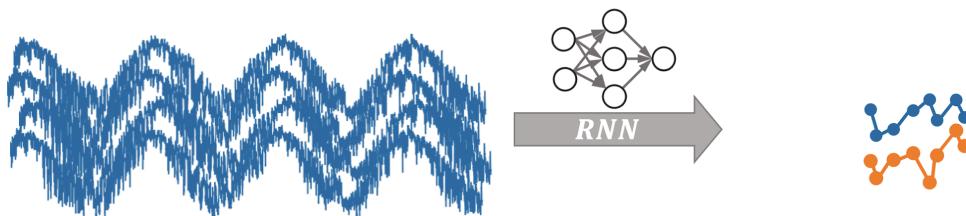
🔮 Prediction with multivariate inputs



👉 Predict `rain` for the `8 next days` based on `multiple features [rain, humidity, temperature, pressure...]`

- `X.shape = (None, 100 days of observation, 4 features)`
- `y.shape = (None, 8 next days, 1 feature = rain)`

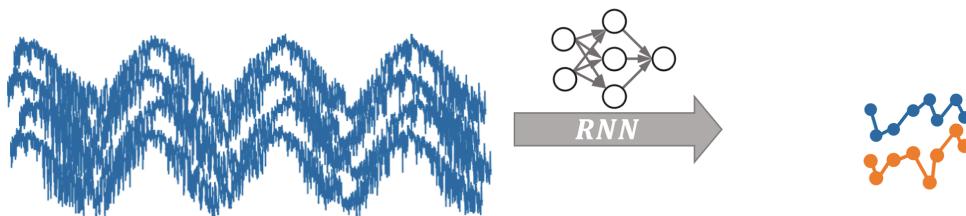
🔮 Predicting multiple outputs



👉 Predict `[rain, temperature]` for the `next 8 days` based on `past [rain, humidity, temperature, pressure...]` features

- `X.shape = (None, 100 days of observation, 4 features)`
- `y.shape = (None, 8 next days, 2 features = [rain, temperature])`

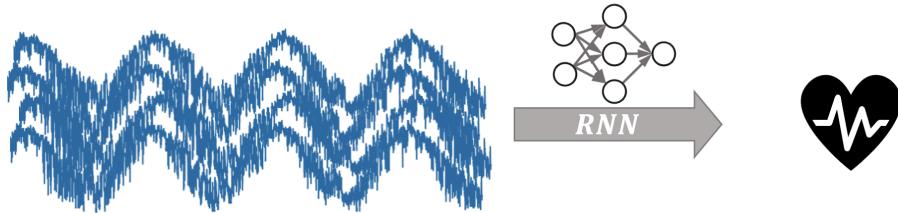
❗ The output is not necessarily similar to the past observations of X!



👉 Predict `[avalanches, accidents]` for the `next 8 days` based on past `[rain, temperature, pressure...]` features

- `X.shape` = (None, 100 days of observation, 4 features = [rain,temperature,pressure])
- `y.shape` = (None, 8 next days, 2 features = [avalanche, accidents])

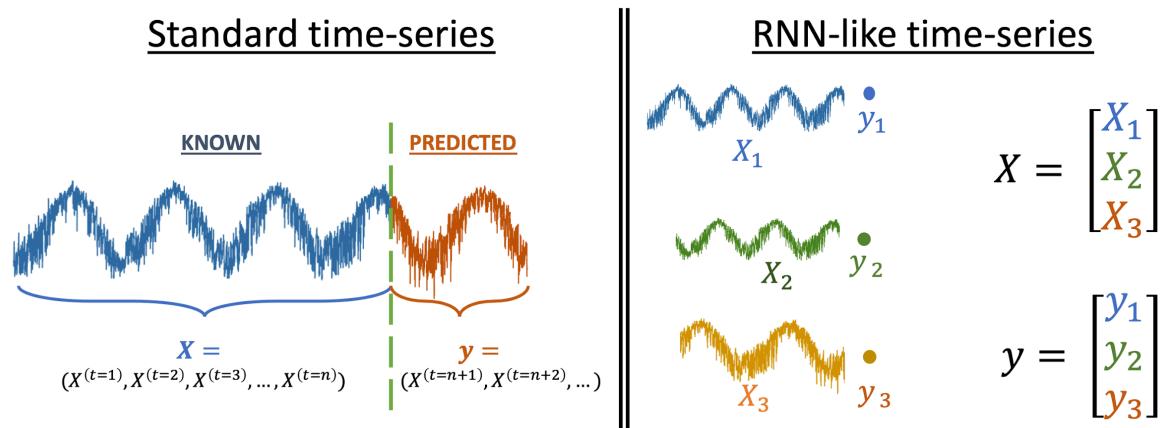
Prediction: using Temporal Data for Classification



👉 Predict the heartbeat category of a patient based on various electrocardiograms (ECG)

- `X.shape` = (None, 100 heartbeats, 4 features)
- `y.shape` = (None, 1 heartbeat, 1 feature = at risk/healthy heart)

Prediction Recap



In Recurrent Neural Networks:

- **Inputs X**
- →
- **sequence of repeated observations** (univariate or multivariate)
 - `n_SEQUENCES` = number of sequences
 - `n_OBSERVATIONS` = `INPUT_LENGTH` = temporal dimension of the input
 - `n_FEATURES` = features you are observing

✓ A dataset in RNN can have **different input lengths**

⚠ We will discuss how we deal with this with padding later in the lecture!

- **Outputs y**
- →

- is whatever you want to predict,
 - **n_SEQUENCES** = number of sequences
 - **OUTPUT_LENGTH** = how many timesteps you want to predict = temporal dimension for the output
 - **n_FEATURES** = features you want to predict

⌚ The feature(s) you want to predict is(are) not naturally similar to the features you are observing.

✖ The shape of what you want to predict must be the same for all the sequences (we'll address this when we talk padding later!)

3.1 Let's code! (Air Pollution Example)

👉 Let's assume that in **different cities**, you observed **weather features** every day during **four days**:

- 🌡️ temperature
- 🌬️ wind speed
- ☁️ air pollution

You are trying to predict the amount of air pollution the **NEXT DAY**.

import numpy as np

--- SEQUENCE A (Paris)

```
day_1 = [10, 25, 60] # OBSERVATION 1 [temperature, speed, pollution]
day_2 = [13, 10, 80] # OBSERVATION 2 [temperature, speed, pollution]
day_3 = [9, 5, 90] # OBSERVATION 3 [temperature, speed, pollution]
day_4 = [7, 0, 100] # OBSERVATION 4 [temperature, speed, pollution]
```

```
sequence_a = [day_1, day_2, day_3, day_4]
```

```
y_a = [110] # Pollution at day 5
```

--- SEQUENCE B (Berlin)

```
sequence_b = [[25, 20, 30], [26, 24, 50], [28, 20, 80], [22, 3, 110]]
y_b = [125]
```

--- SEQUENCE C (London)

```
sequence_c = [[15, 10, 60], [25, 20, 65], [35, 10, 75], [36, 15, 70]]
y_c = [75]
```

```
X = np.array([sequence_a, sequence_b, sequence_c]).astype(np.float32)
y = np.expand_dims(np.array([y_a, y_b, y_c]).astype(np.float32), axis=-1)
```

```
print(X.shape)
print(y.shape)
```

(3, 4, 3)
(3, 1, 1)

```
# 0- Imports
from tensorflow.keras import Sequential, Input, layers
from tensorflow.keras.optimizers import Adam

# 1- RNN Architecture
model = Sequential()
model.add(Input(shape=(4, 3)))
model.add(layers.SimpleRNN(units=2, activation='tanh'))
model.add(layers.Dense(1, activation="linear"))

# 2- Compilation
model.compile(loss='mse',
               optimizer=Adam(learning_rate=0.5)) # very high lr so we can converge with such a small
dataset

# 3- Fit
model.fit(X, y, epochs=10000, verbose=0)

# 4- Predict
model.predict(X) # One prediction per city
```

1/1 ————— 0s 89ms/step

```
array([[117.499985],
       [117.499985],
       [ 75.000015]], dtype=float32)
```

3 RNNs Under the Hood

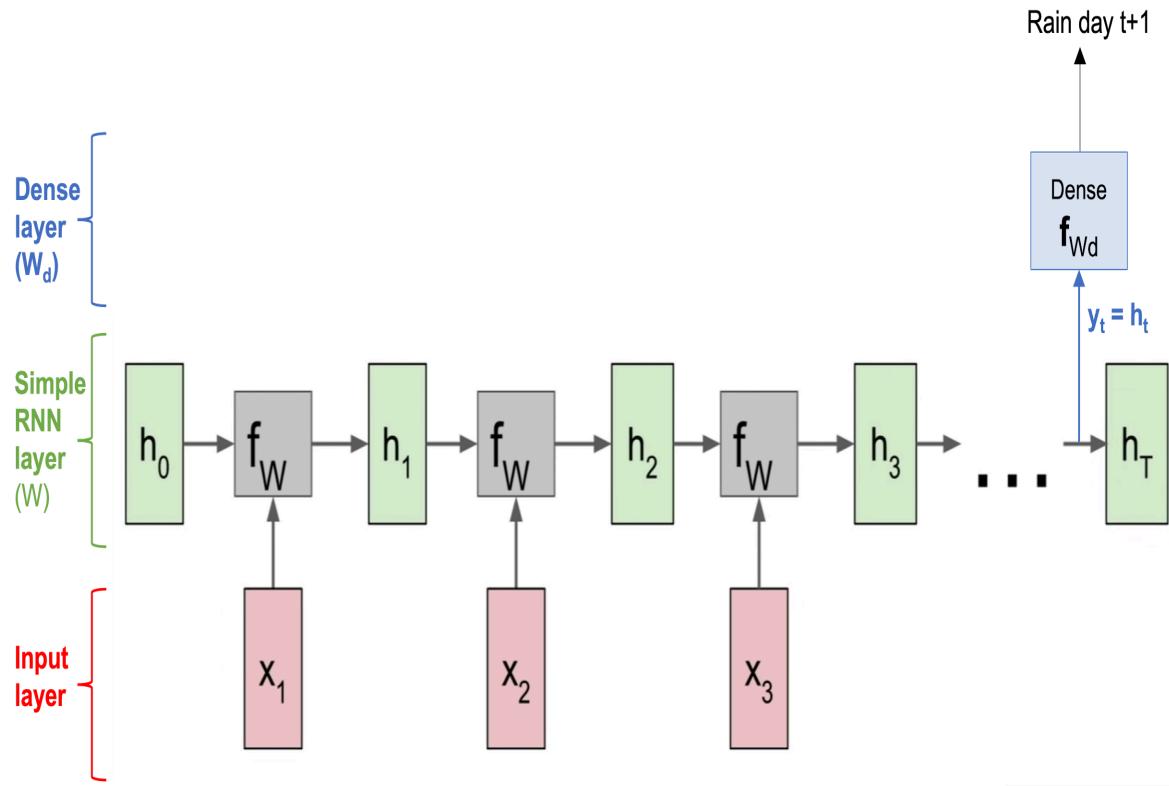
The RNN is fed **one observation**

X
(
k
)
at a time (forward in time).

It maintains an **internal state**

h

that is updated at each time step.

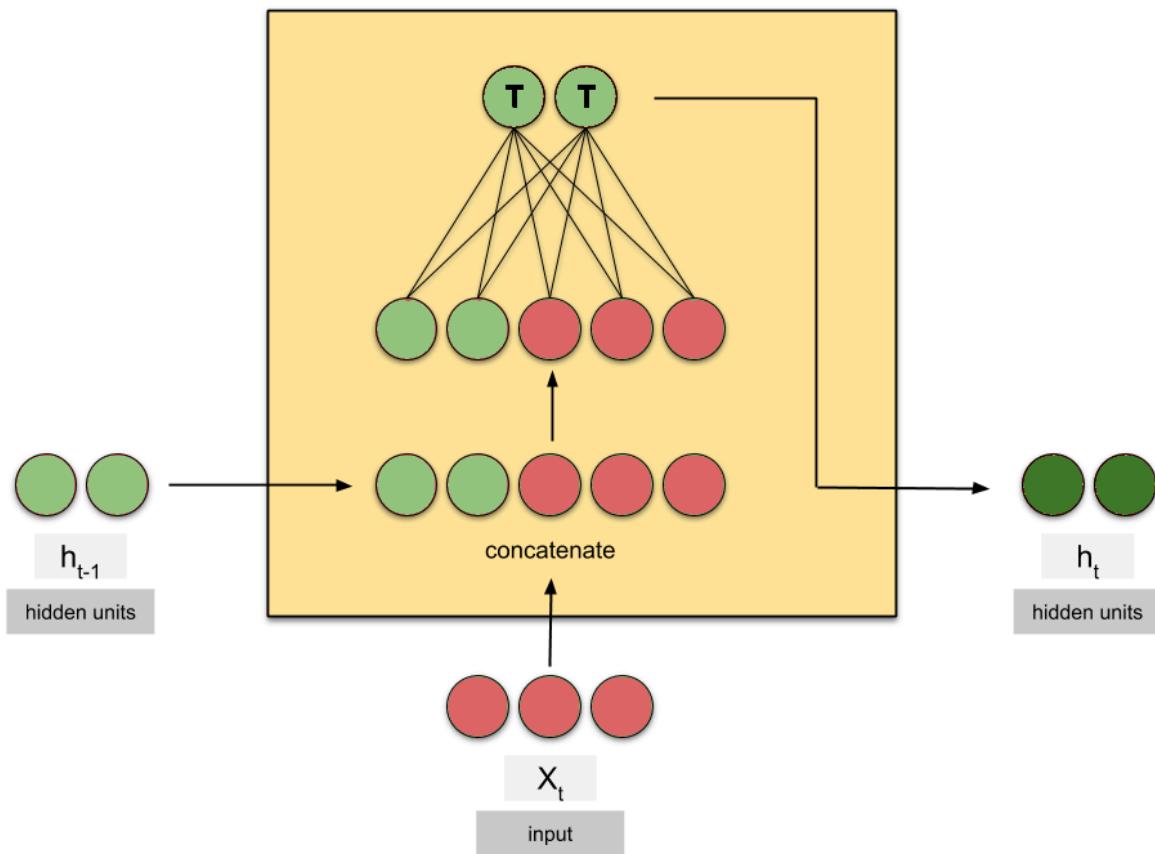


Let's zoom inside this function

f

W

for one time step



- The **current observation**
- **x**
- **(**
- **t**
- **)**
- and the **previous internal state**
- **h**
- **(**
- **t**
- **-**
- **1**
- **)**
- are **concatenated** and fed into a **weights matrix**...
- ...to output the **new internal state**
- **h**
- **(**
- **t**
- **)**
- **!**

```
model.summary()
```

```
Model: "sequential"
```

#	Layer (type)	Output Shape	Param
	simple_rnn (SimpleRNN)	(None, 2)	12
	dense (Dense)	(None, 1)	3

```
Total params: 47 (192.00 B)
```

```
Trainable params: 15 (60.00 B)
```

```
Non-trainable params: 0 (0.00 B)
```

```
Optimizer params: 32 (132.00 B)
```

Recall how we defined our SimpleRNN layer as:
model.add(layers.SimpleRNN(units=2, activation='tanh'))

One single RNN layer has a matrix

W

of trainable params, such that

h

(

t

)

=

f

w

```

( h
( t
-
1
)
,
x
(
t
)
)

```

- The **same** weights
 - W
 - are used at every time step...
 - ...to **recursively** compute its **internal state**
 - h
 -
-

We can *count* the number of weights of the dense network above as

$$n_h(n_h + n_x)$$

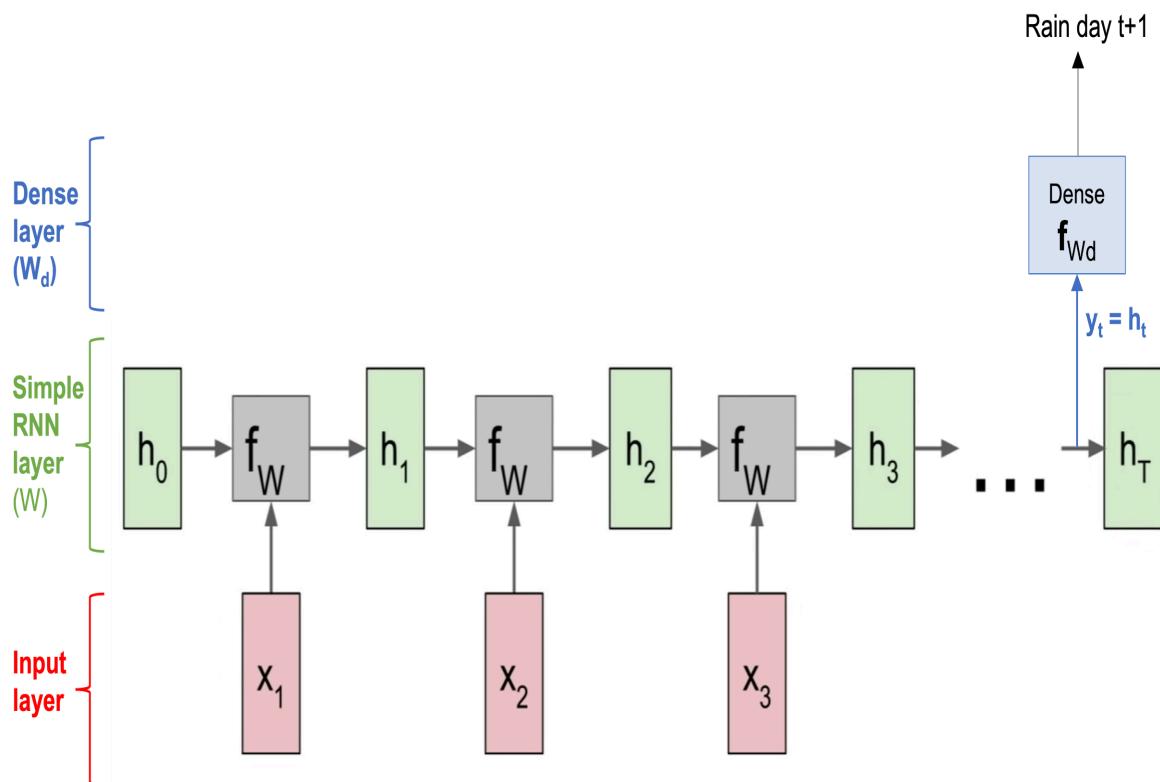
+

1

)

.

- n
- h
- = number of **units** in our RNN layer
- n
- x
- = number of **features** measured per time-step (e.g. [temperature, humidity, ...])



! Two important notes! !

1. In a similar way to the number of weight remaining the same in CNNs regardless of image size, the number of weights in a Recurrent Layer has nothing to do with the length of our sequences (4 days in our example).
2. ~~X~~
3. y
4. $($
5. t
6. $)$

7. is NOT the prediction/ the target.
8. **y**
9. **(**
10. **t**
11. **)**
12. is a vector of size **RNN_units** (
13. **n**
14. **h**
15. **=**
16. **2**
17. in our example)...
 ... used as an input to the Dense layer to compute the **rain** at time
18. **t**
19. **+**
20. **1**

How can we improve our model?

Increase the number of RNN units

h

?

From `model.add(layers.SimpleRNN(units=2, activation='tanh'))`

to `model.add(layers.SimpleRNN(units=10, activation='tanh'))`

👉 **Intuition:** With ****10**** RNN units rather than just 2, this model:

- will try to *capture **10** interesting temporal features* from the time-series
 - (maybe: *mean, rate of increase, complex auto-regressive feature, etc...who knows!*)
- and combine them into 1 value for our regression task

💡 The number of units can be seen as the **number of memories about features maintained in parallel.**

Add another hidden layer?

From:

```
model.add(layers.SimpleRNN(units=2, activation='tanh'))
model.add(layers.Dense(1, activation="linear"))
```

To:

```
model.add(layers.SimpleRNN(units=2, activation='tanh'))
```

```
model.add(layers.Dense(10, activation="linear"))
model.add(layers.Dense(1, activation="linear"))
```

Add another RNN layer?

Much like chaining together Convolutional Layers for CNNs, we can chain together RNN layers 💪

From:

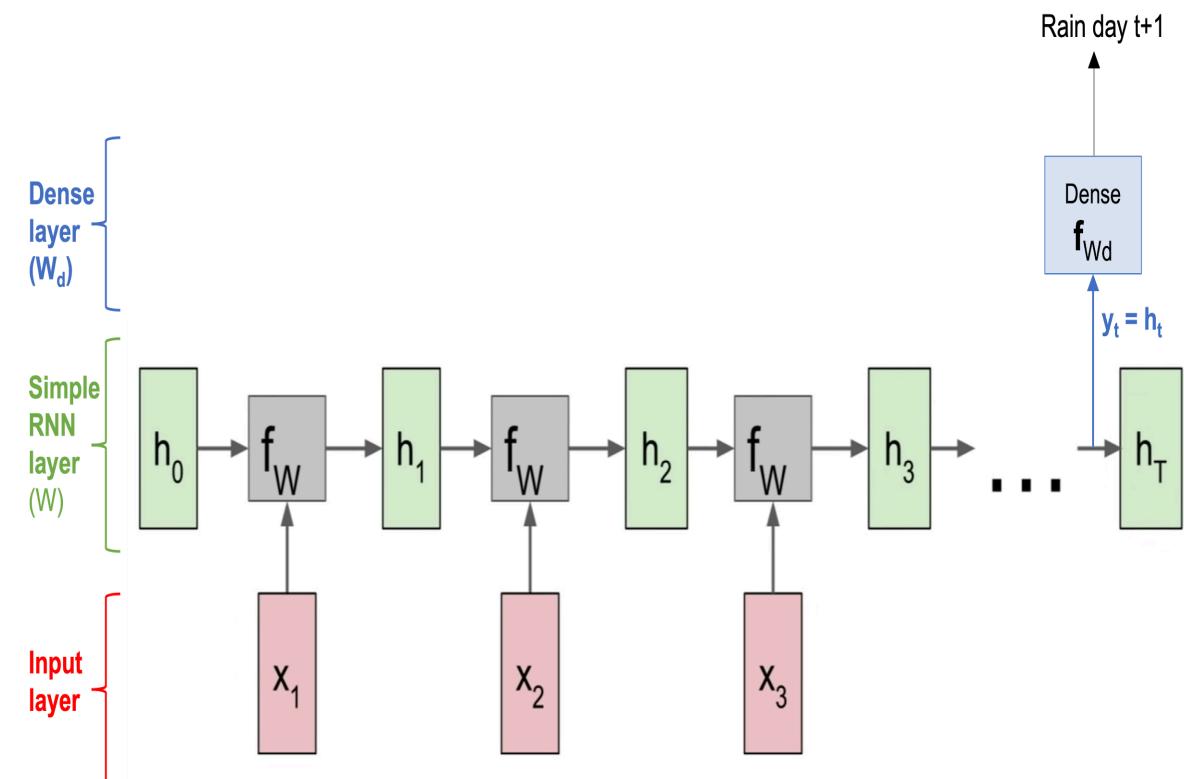
```
model.add(layers.SimpleRNN(units=2, activation='tanh'))
model.add(layers.Dense(1, activation="linear"))
```

To:

```
model.add(layers.SimpleRNN(units=2, return_sequences=True,
activation='tanh'))
model.add(layers.SimpleRNN(units=2, activation='tanh'))
model.add(layers.Dense(1, activation="linear"))
```

! Take note! ! We've added `return_sequences=True`; but why?

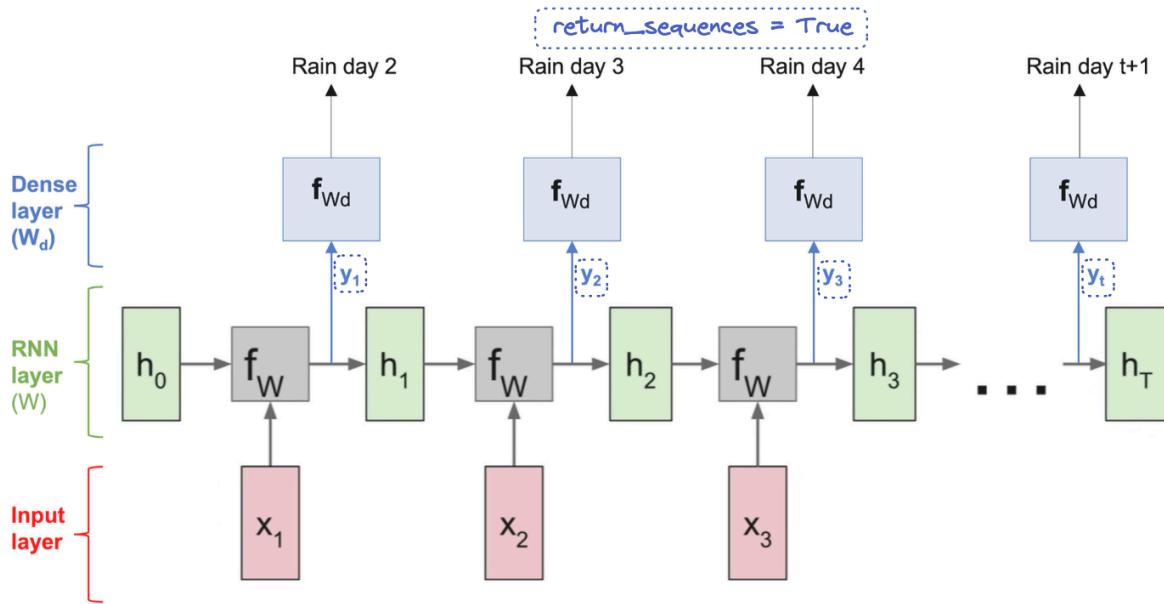
Our initial approach:



`return_sequences=True` outputs

```
y
(
t
)
```

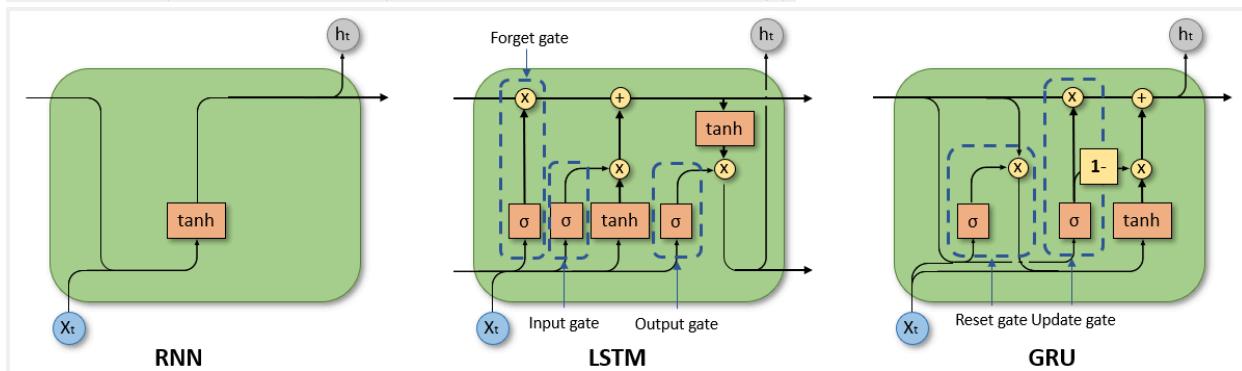
at each time step.



These can be fed into another RNN layer!

Change the type of RNN we have in our model!

```
model.add(layers.LSTM(units=5, activation='tanh', input_shape=(4, 3)))
model.add(layers.Dense(10, activation="linear"))
model.add(layers.Dense(1, activation="linear"))
```



```
model = Sequential()
model.add(Input(shape=(4, 3)))
model.add(layers.LSTM(units=2, activation='tanh'))
model.add(layers.Dense(1, activation="linear"))
model.summary()
```

```
Model: "sequential_1"
```

Layer (type)	Output Shape	Param
#		
lstm (LSTM)	(None, 2)	48

dense_1 (Dense)	(None, 1)	3
-----------------	-----------	---

Total params: 51 (204.00 B)

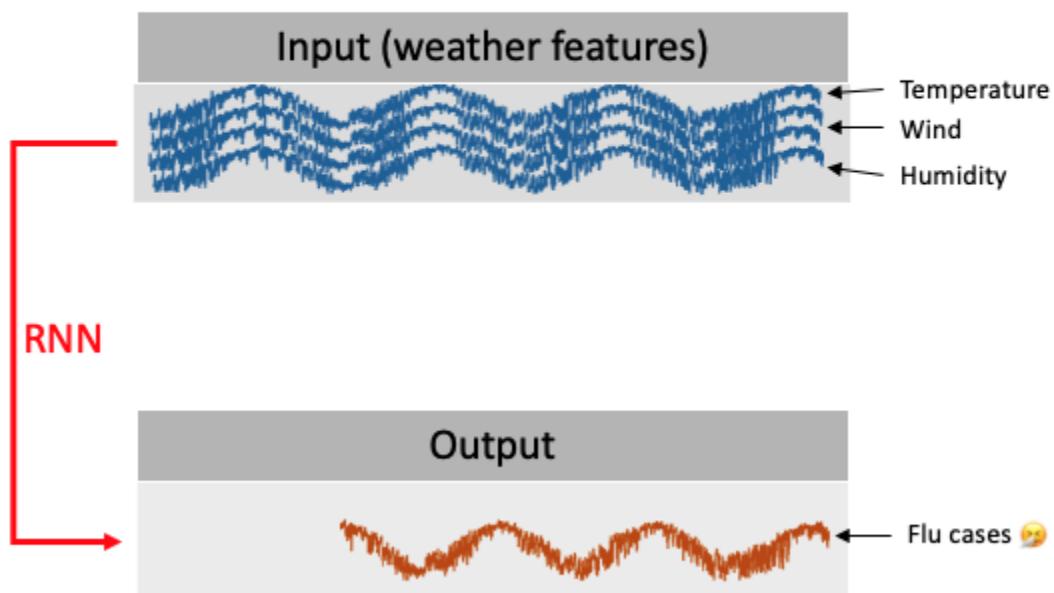
Trainable params: 51 (204.00 B)

Non-trainable params: 0 (0.00 B)

⚠️ N.B. We now have 4 times as many weights in our RNN layer thanks to the additional gates! ⚠️

Predicting an entire sequence

Sequence A (Paris)



For each sequence/city:

X

(

t

)

= weather observed at day

t

y

(

t

)

= Flu cases reported at day

t

Why not use

y

at previous time steps as a feature?

! To predict a sequence, `y_train` need to be sequences too !
`import numpy as np`

--- SEQUENCE A (Paris)

```
sequence_a = [[10, 25, 50], # OBS day 1  
              [13, 10, 70], # OBS day 2  
              [ 9,  5, 90], # OBS day 3  
              [ 7,  0, 95]] # OBS day 4
```

```
y_a = [70,  # flu cases day 1  
       90,  # flu cases day 2  
       95,  # flu cases day 3  
      110,] # flu cases day 4
```

```

# --- SEQUENCE B (Berlin)
sequence_b = [[25, 20, 30], [26, 24, 50], [28, 20, 80], [22, 3, 110]]
y_b = [50, 80, 110, 125]

# --- SEQUENCE C (London)
sequence_c = [[15, 10, 60], [25, 20, 65], [35, 10, 75], [36, 15, 70]]
y_c = [65, 75, 70, 30]

X = np.array([sequence_a, sequence_b, sequence_c]).astype(np.float32)
y = np.expand_dims(np.array([y_a, y_b, y_c]).astype(np.float32), axis=-1)

print(X.shape)
print(y.shape)

```

(3, 4, 3)

(3, 4, 1)

! Remember the wording !

```

X.shape = (n_SEQUENCES, n_OBSERVATIONS, n_FEATURES)
X.shape = (3 cities, 4 days, 3 features)
y.shape = (3 cities, 4 days, 1 feature predicted each day)
model_2 = Sequential()

model_2.add(layers.SimpleRNN(2, return_sequences=True))
model_2.add(layers.Dense(1, activation='relu'))

model_2.compile(loss='mse', optimizer='adam')
model_2.fit(X, y, verbose=0)

```

<keras.src.callbacks.history.History at 0x304fe4740>

print("y_pred shape:", model_2.predict(X).shape)

1/1 ━━━━━━━━━━ 0s 163ms/step

y_pred shape: (3, 4, 1)

4 NLP: Why we need RNNs

Because order matters!

An obvious example:

"The teacher inspired the students with a passionate speech about deep learning"

"The students inspired the teacher with a passionate speech about deep learning"

A more subtle example:

"The company hired talented employees and achieved great success."

"The company achieved great success and hired talented employees."

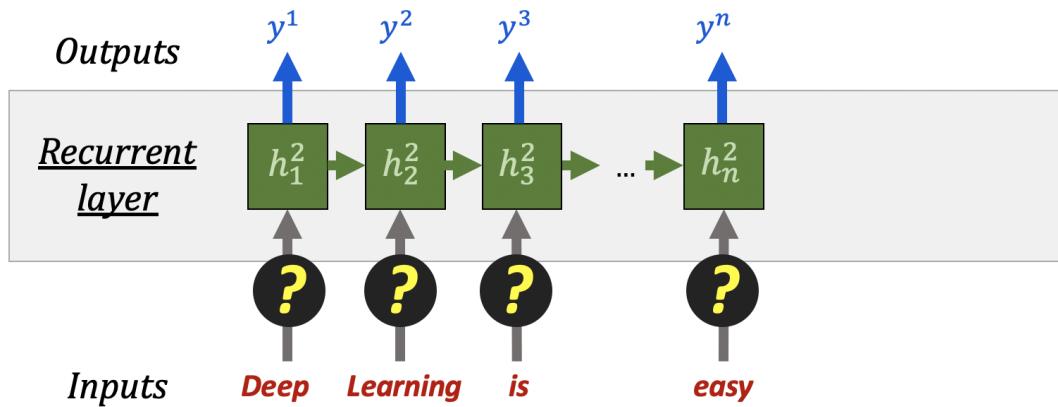
Let's give ourselves a task to solve by the end of the lecture!

Our X will be sentences and our y will be whether or not those sentences are positive or negative!

How do we go about this?

"This movie is the best thing I've seen in my life" -> "I had to leave this film because I was so bored" -> "Anything with Nicolas Cage in it is a masterpiece: 5/5" ->

How to feed Recurrent Neural Networks with words?



We have to convert words into numbers somehow!

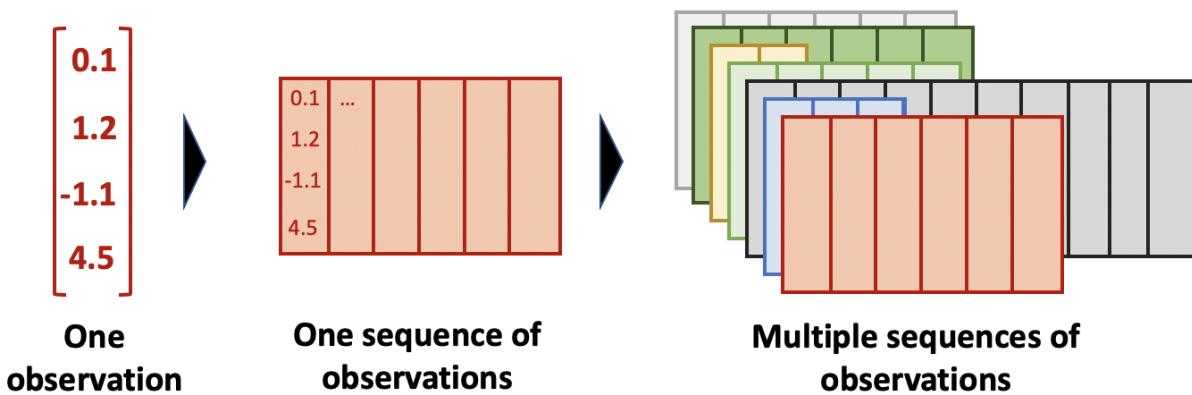
Text is also a form of recurrent data, where a sentence is a sequence of words, each being the *observation*. For now, let's imagine we can express each **word** as **four numbers** (or a point/ vector in 4-D space)

$$\begin{array}{c}
 \begin{bmatrix} 0.1 \\ 1.2 \\ -1.1 \\ 4.5 \end{bmatrix} \quad \begin{bmatrix} 0.3 \\ 2.2 \\ -0.7 \\ 2.9 \end{bmatrix} \quad \begin{bmatrix} 0.0 \\ 1.8 \\ -0.2 \\ 1.9 \end{bmatrix} \quad \cdots \quad \begin{bmatrix} 1.0 \\ -2.8 \\ 1.9 \\ 0.3 \end{bmatrix} \\
 \text{Observation at time 1} \quad \text{Observation at time 2} \quad \text{Observation at time 3} \quad \quad \quad \text{Observation at time } T
 \end{array}$$

Deep Learning is ... easy

Word 1 Word 2 Word 3 ... Word N

What will our X look like?



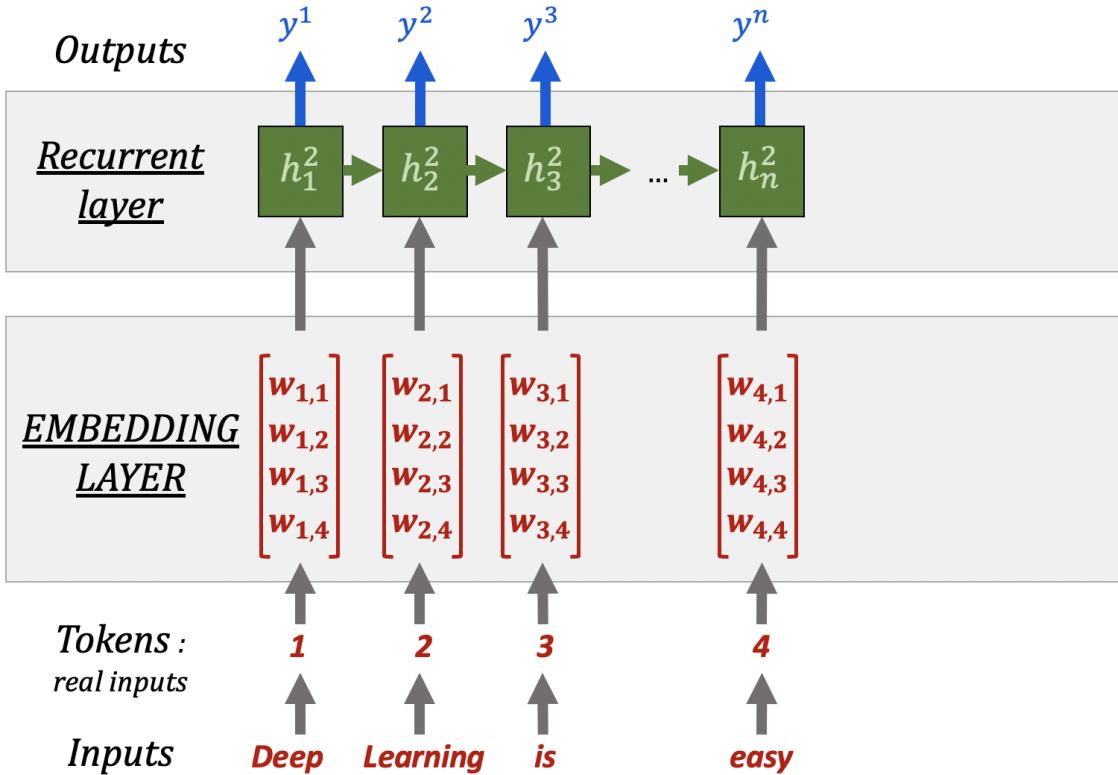
X is just a sequence of observations! Univariate, multivariate, an image, ... The point is that the observation is repeated through time!

```
X.shape = (N_SEQUENCES, N_OBSERVATIONS, N_FEATURES)
```

Warning: The number of observations can vary from one sequence to another! (This is when you pad - we'll get back to that)

First, we'll need tokens

Neural networks' can't work directly on words, so you will still need to provide it with tokens.



```
import numpy as np
from tensorflow.keras.preprocessing.sequence import pad_sequences
from tensorflow.keras.preprocessing.text import Tokenizer
```

```
### Let's create some mock data
def get_mock_up_data():
    sentence_1 = 'This movie was awful'
    sentence_2 = 'I loved every moment of this movie!'
    sentence_3 = 'I want a refund; I was so bored!'
```

```
X = [sentence_1, sentence_2, sentence_3]
y = np.array([0., 1., 0.])
```

```
### Let's tokenize the vocabulary
tk = Tokenizer()
tk.fit_on_texts(X)
vocab_size = len(tk.word_index)
print(f'There are {vocab_size} different words in your corpus')
X_token = tk.texts_to_sequences(X)
```

```
### Pad the inputs
X_pad = pad_sequences(X_token, dtype='float32', padding='pre')

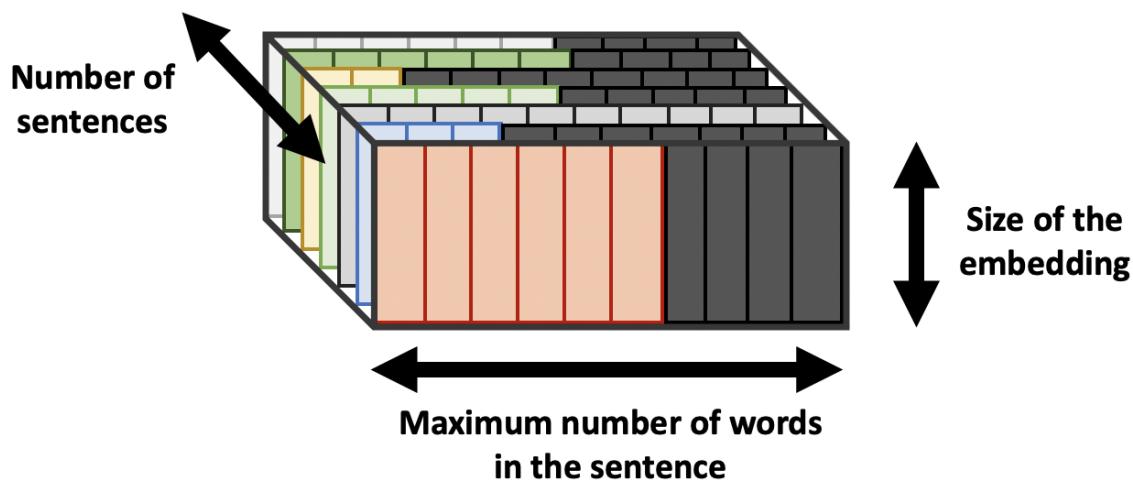
return X_pad, y, vocab_size
```

```
X_pad, y, vocab_size = get_mock_up_data()  
print("X_pad.shape", X_pad.shape)  
X_pad
```

There are 14 different words in your corpus
X_pad.shape (3, 8)

```
array([[ 0.,  0.,  0.,  0.,  2.,  3.,  4.,  5.],  
       [ 0.,  1.,  6.,  7.,  8.,  9.,  2.,  3.],  
       [ 1., 10., 11., 12.,  1.,  4., 13., 14.]], dtype=float32)
```

What's going on with the `pad_sequences` function?



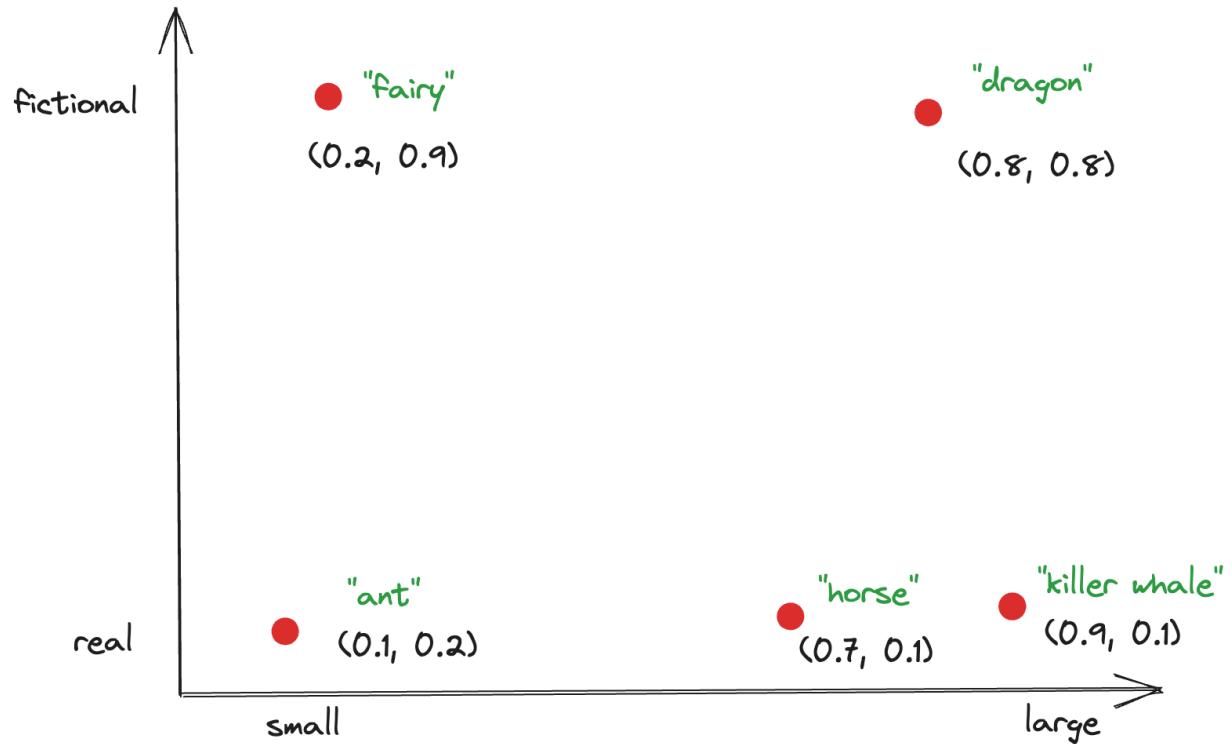
🤔 To "post" or to "pre"? The eternal question: there's no definitive answer and it varies by architecture, but results seem to lean towards using "["pre"](#)"

Now we need to make the leap from tokens to vectors

We want each word to be represented by a vector ↗ of *chosen* length. How?

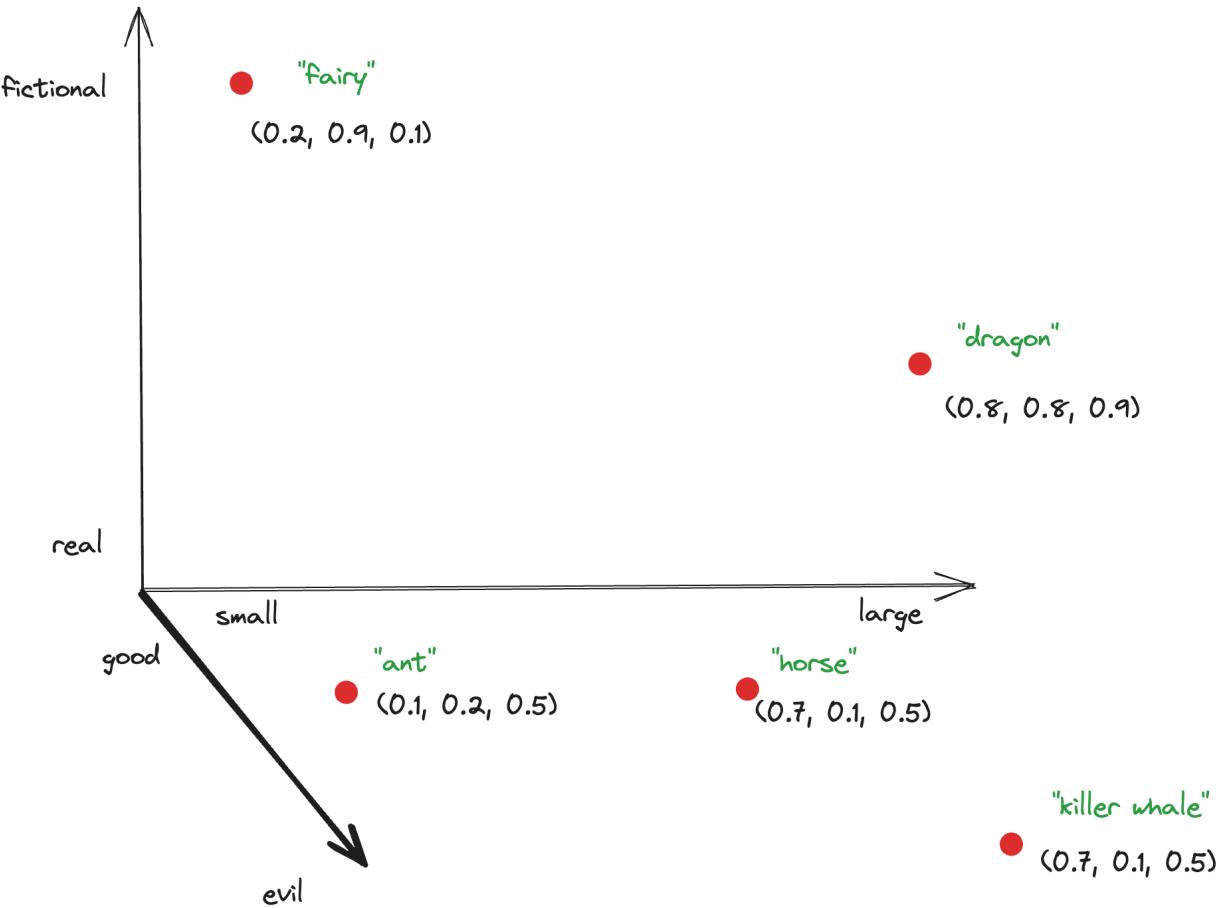
Consider a 2D embedding

To further specify the meaning of each word, we add a dimension: we rank them according to their relative abstraction.



We can't differentiate much between "horse" and "killer whale" 😞

Adding a 3rd dimension



Arithmetic on words?

If we embed correctly, we could dream of performing mathematical operations on the embeddings, which would have meaning in terms of natural language.

E.g. we could express the sentence "Queen is to king as man is to..." as

V

(

Q

u

e

e

n

)

-

v

(

K

i

n

g

)

+

v

(

M

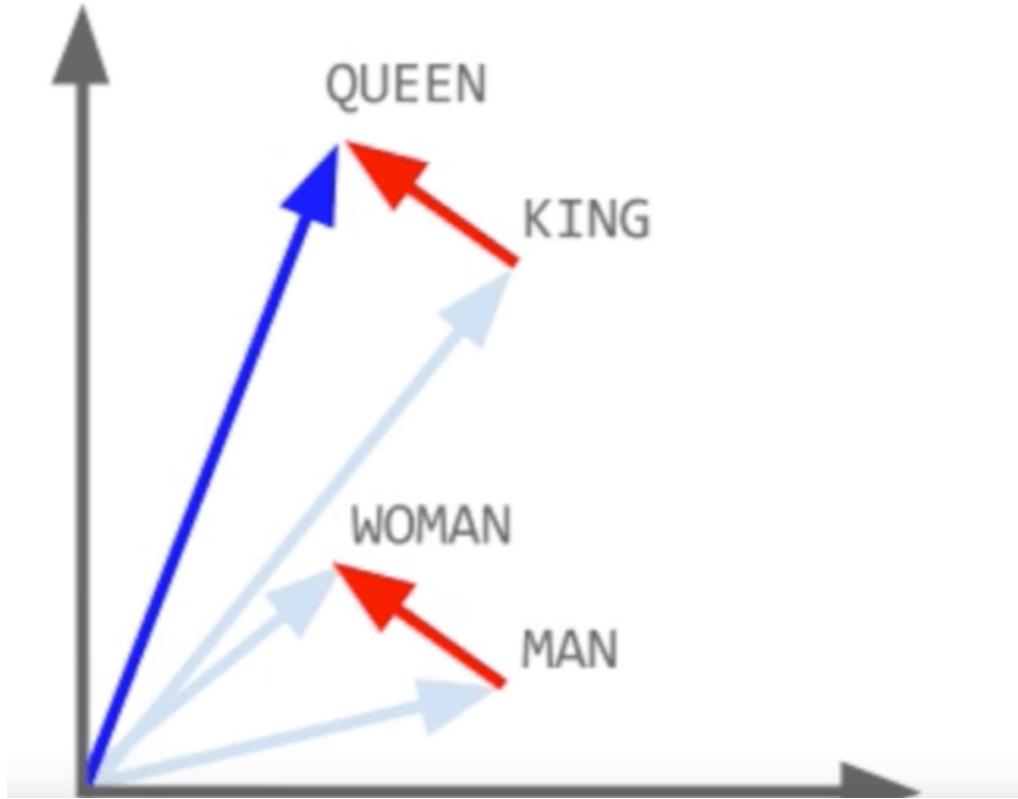
a

n

)

=

?



Here, as the red vectors are similar, we could expect

V
(
Q
u
e
e
n
)
-
V
(

K

i

n

g

)

=

V

(

W

o

m

a

n

)

-

V

(

M

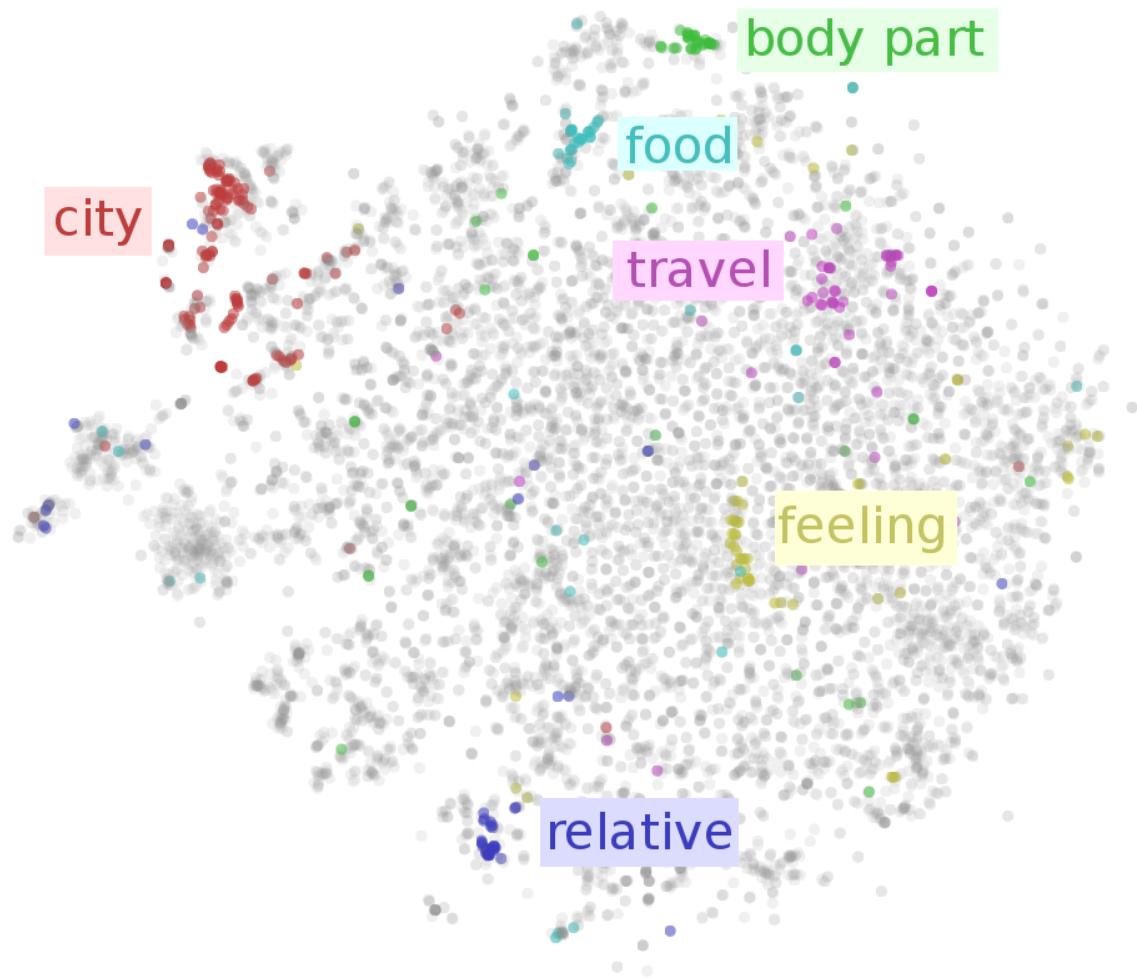
a

n

)

So the answer we're looking for might be **Woman**.

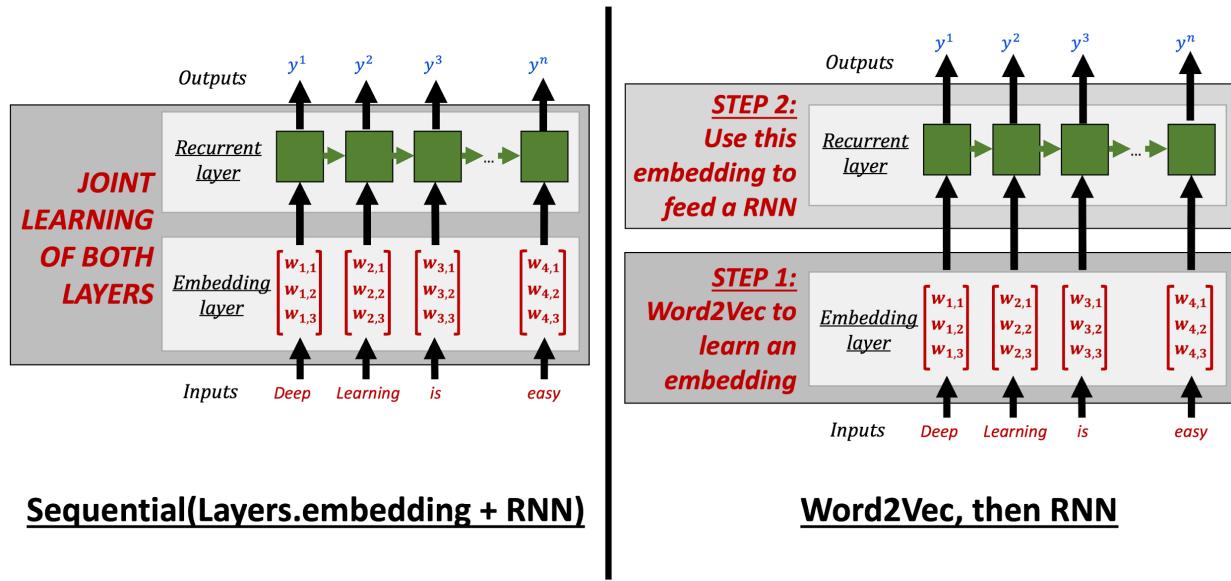
What makes a good embedding?



Semantically *close* words are mathematically *close* in this space

❗ Usually, word embedding spaces are from **30** up to **300** dimensions ❗
🤔 But wait ... how do you choose these numbers? Surely we're not scoring all of our words on n difference axes?

Two options: learn with `layers.Embedding` or with Word2Vec



👉 The left option allows you to have a representation that is perfectly suited to your task! However, it increases the number of parameters to learn, and thus:

- the time of each epoch (more parameters to optimize during back-propagation)
- the time to converge (because more parameters to find overall)

👉 On the other hand, word2vec is an unsupervised learning method not specifically designed for your task (may be sub-optimal) but training it is very fast! You will also be able to optimize your RNN faster as you'll have less parameters.

- ! Prefer Word2Vec on small corpus (esp. with transfer learning)

First option: layers.Embedding

```
Embedding(input_dim=VOCAB_SIZE,
          input_length=MAX_SENTENCE_LENGTH,
          output_dim=EMBED_DIM,
          mask_zero=True)
```

We can insert this layer after our tokenization and all of our words are embedded as the model trains!

Second option: Word2Vec

Source Text	Training Samples
The quick brown fox jumps over the lazy dog. →	(the, quick) (the, brown)
The quick brown fox jumps over the lazy dog. →	(quick, the) (quick, brown) (quick, fox)
The quick brown fox jumps over the lazy dog. →	(brown, the) (brown, quick) (brown, fox) (brown, jumps)
The quick brown fox jumps over the lazy dog. →	(fox, quick) (fox, brown) (fox, jumps) (fox, over)

💡 We use unsupervised learning to look at the words around each chosen word and assume that these will be relevant to the word we're interested in!

Instead of training it on your training set (especially if it is very small), you can directly load a pretrained embedding

`import gensim.downloader`

```
print(list(gensim.downloader.info()['models'].keys()))
```

```
model_wiki = gensim.downloader.load('glove-wiki-gigaword-50')
```

```
['fasttext-wiki-news-subwords-300', 'conceptnet-numberbatch-17-06-300',  
'word2vec-ruscorpora-300', 'word2vec-google-news-300', 'glove-wiki-gigaword-50',  
'glove-wiki-gigaword-100', 'glove-wiki-gigaword-200', 'glove-wiki-gigaword-300', 'glove-twitter-25',  
'glove-twitter-50', 'glove-twitter-100', 'glove-twitter-200', '__testing_word2vec-matrix-synopsis']
```

We can even try out our aforementioned example with code!

```
# N.B. Words not in glove-wiki-gigaword-50 will not have vectors computed  
example_1 = model_wiki["queen"] - model_wiki["king"] + model_wiki["man"]
```

```
model_wiki.most_similar(example_1)
```

```
[('woman', 0.8903914093971252),  
 ('girl', 0.8453726768493652),  
 ('man', 0.8301756381988525),  
 ('her', 0.7845831513404846),  
 ('boy', 0.7763065099716187),
```

```
('she', 0.7619765400886536),  
('herself', 0.7597628831863403),  
('blind', 0.7296754717826843),  
('mother', 0.7230340242385864),  
('blonde', 0.713614284992218)]
```

What about "good is to evil as cold is to..."?

```
example_2 = model_wiki["good"] - model_wiki["evil"] + model_wiki["cold"]
```

```
model_wiki.most_similar(example_2)
```

```
[('warm', 0.7870427966117859),  
('dry', 0.7643216848373413),  
('hot', 0.750431478023529),  
('cool', 0.7491167187690735),  
('weather', 0.7476345896720886),  
('getting', 0.7447267770767212),  
('good', 0.7442173361778259),  
('cold', 0.7425146102905273),  
('low', 0.7231549024581909),  
('little', 0.7223491668701172)]
```

5 Coded example: Sentiment analysis

First let's load up our data 🍔

```
import tensorflow as tf  
from tensorflow import keras  
from tensorflow.keras.datasets import imdb  
from tensorflow.keras.preprocessing import sequence  
from tensorflow.keras import Sequential  
from tensorflow.keras.layers import Embedding, LSTM, Dense  
  
# Set parameters  
max_features = 10000 # Maximum number of words to get out of our imdb data  
max_len = 500 # Maximum sequence length  
embedding_dim = 50 # Dimensionality of word embeddings  
  
# Load the IMDb dataset  
(x_train, y_train), (x_test, y_test) = imdb.load_data(num_words=max_features)
```

```
Downloading data from https://storage.googleapis.com/tensorflow/tf-keras-datasets/imdb.npz  
17464789/17464789 ━━━━━━━━━━━━━━ 2s 0us/step
```

Next we need to pad 📄

```
# Pad sequences to a fixed length
```

```
x_train = sequence.pad_sequences(x_train, maxlen=max_len)
```

```
x_test = sequence.pad_sequences(x_test, maxlen=max_len)
```

Now we'll construct and compile our model 

Build the model

```
model = Sequential()  
model.add(Embedding(max_features, embedding_dim))  
model.add(LSTM(64))  
model.add(Dense(1, activation='sigmoid'))
```

Compile the model

```
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
```

Finally we fit and evaluate 

Train the model

```
model.fit(x_train, y_train, batch_size=128, epochs=3, validation_data=(x_test, y_test))
```

Evaluate the model

```
loss, accuracy = model.evaluate(x_test, y_test)  
print(f'Test loss: {loss:.4f}')  
print(f'Test accuracy: {accuracy:.4f}')
```

Epoch 1/3

196/196  **75s** 380ms/step - accuracy: 0.6632 - loss: 0.5910 - val_accuracy: 0.8569 - val_loss: 0.3400

Epoch 2/3

196/196  **76s** 387ms/step - accuracy: 0.8904 - loss: 0.2785 - val_accuracy: 0.8722 - val_loss: 0.3052

Epoch 3/3

196/196  **75s** 382ms/step - accuracy: 0.9276 - loss: 0.2000 - val_accuracy: 0.8694 - val_loss: 0.3089

782/782  **21s** 27ms/step - accuracy: 0.8683 - loss: 0.3109

Test loss: 0.3089

Test accuracy: 0.8694

Done! We can predict sentiment with a very high degree of accuracy and minimal pre-processing!

6 Further reading

-  [Stanford Cheat Sheet](#)
-  [Medium - Illustrated Guide to RNN](#)
-  [Medium - Illustrated Guide to LSTM and GRUs](#)
-  [RNN explained with 3*3 matrices - 21min](#)
-  [RNN & LSTM explained without math - 24min](#)

Deep Learning courses - videos

- [CS224N: Natural Language Processing with Deep Learning](#) by Standford University (recommended)
- [92 step-by-step videos](#) by Hugo Larochelle.
- [Deep Learning course](#), by Andrew Ng.
- [MIT Introduction to Deep Learning](#), from MIT.
- [Deep Learning for Computer Vision](#), from Michigan University.

Deep Learning courses - books

- [Deep Learning book](#), by Ian Goodfellow, Yoshua Bengio, Aaron Courville
- [Deep Learning with Python](#), by François Chollet231 Lecture on RNN - 1h20](<https://www.youtube.com/watch?v=6niqTuYFZLQ>)