

## DIFFERENCE BETWEEN CONVOLUTION LAYER AND FULLY DENSELY CONNECTED LAYER

Dense layers learn "global" patterns in their input feature space (e.g. for a MNIST digit, patterns involving all pixels)

Conv layers learn "local" patterns (e.g. for a MNIST digit, patterns involving small 2D windows of inputs)

## KEY CHARACTERISTICS/PROPERTIES OF CONVNETS:

- Convnets learn patterns that are 'translation invariant':  
i.e. After learning a pattern in the lower right corner of a picture, a convnet can recognise it anywhere

- Convnets learn 'spatial hierarchies' of patterns:

A first convolution layer will learn small local patterns such as edges, a second conv. layer will learn larger patterns made of features of the first layer and so on. This allows convnets to learn increasingly complex visual concepts

## CONVOLUTION LAYER

- Convolution layer is defined by two key params:

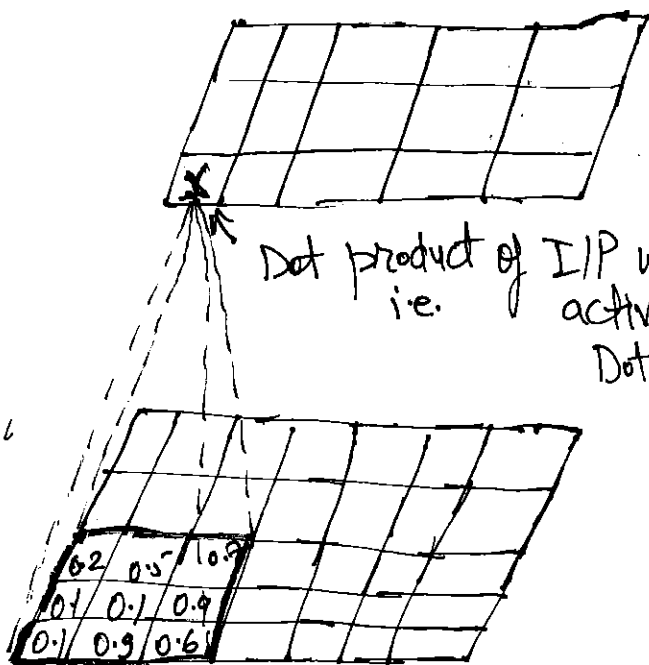
i) Size of patches extracted from inputs = filter size  
(Typically  $3 \times 3$ )

ii) # of filters = # of o/p feature maps  
conv. kernels

Note: If the I/P to conv. layer is n-dimensional (e.g. 3),  
the filter depth is n-dimensional

Example: in case of RGB image, the first conv  
layer's filters will have depth 3

And if conv2D layer is used then  
o/p of this layer is 2D and the # of  
filters determine the depth of o/p from this  
layer i.e. # of feature maps



Dot product of I/P with filter followed by activation  
i.e. activation (dot product I/P with filter)

Dot product results in scalar hence one value

0.1	0.9	0
0.0	0.9	0.1
0.0	0.7	0

→ filter ( $3 \times 3$ )

Stride: → The way the filter moves from one position to next on the I/P

→ The factor by which width and ht. of I/P feature maps are downsampled (in addition to any changes induced by padding)

e.g. stride = 2 means width and ht. of I/P feature map is downsampled by factor of 2

→ Typically we keep stride = 1 in conv. layer

Padding: → When the filter moves across the I/P feature map, the output feature map's width & ht. are still reduced (even when stride = 1)

a	b	c
d	e	f
g	h	i

I/P feature map  
3x3

using 2x2  
filters with  
stride = 1

ab	bc
de	ef
gh	hi

2x2  
O/P feature map

Zero Padding → To make O/P feature map's width & ht = I/P feature map's width & ht  
i.e. I/P feature map is padded with cells left, right, top & bottom in such a way that O/P feature map's width & ht = I/P feature map's width & ht

Note: Depth of O/P feature map depends on # of filters

→ Filters = receptive field

e.g. if the filter is

0	1	0
0	1	0
0	1	0

will detect

central vertical line  
(if the image is grayscale  
i.e. matrix of values b/w  
0 & 1)

Thus a layer full of neurons using the same filter  
O/P a feature map which highlights the areas  
in an image that activate the filter the most.

→ The values in the filter are the wts. that are  
learnt during the training process.

→ one filter results in one feature map  
That's why all neurons in the feature map share  
the same parameters

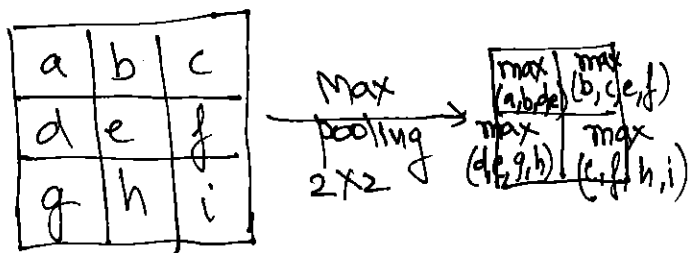
→ When one conv. layer is directly connected to another  
conv. layer (i.e. without any other layer in between),  
the O/P feature map of one conv. layer goes as  
I/P feature map to next conv. layer

## POOLING LAYER

→ Pooling layer is defined by:

- i) size of window/patches. (typically  $2 \times 2$ )
- ii) stride (typically 2)

→ Pooling consists of extracting windows from I/P feature map and performing the pooling operation (max, avg) on the extracted window  $\leftarrow$  o/p of each channel (feature map)



→ Downsamples width & ht of I/P feature map by the factor of stride, does not downsample depth/# of feature maps (typically)

→ However, downsampling by depth can also be done (not provided in Keras, need to use low-level TF APIs)

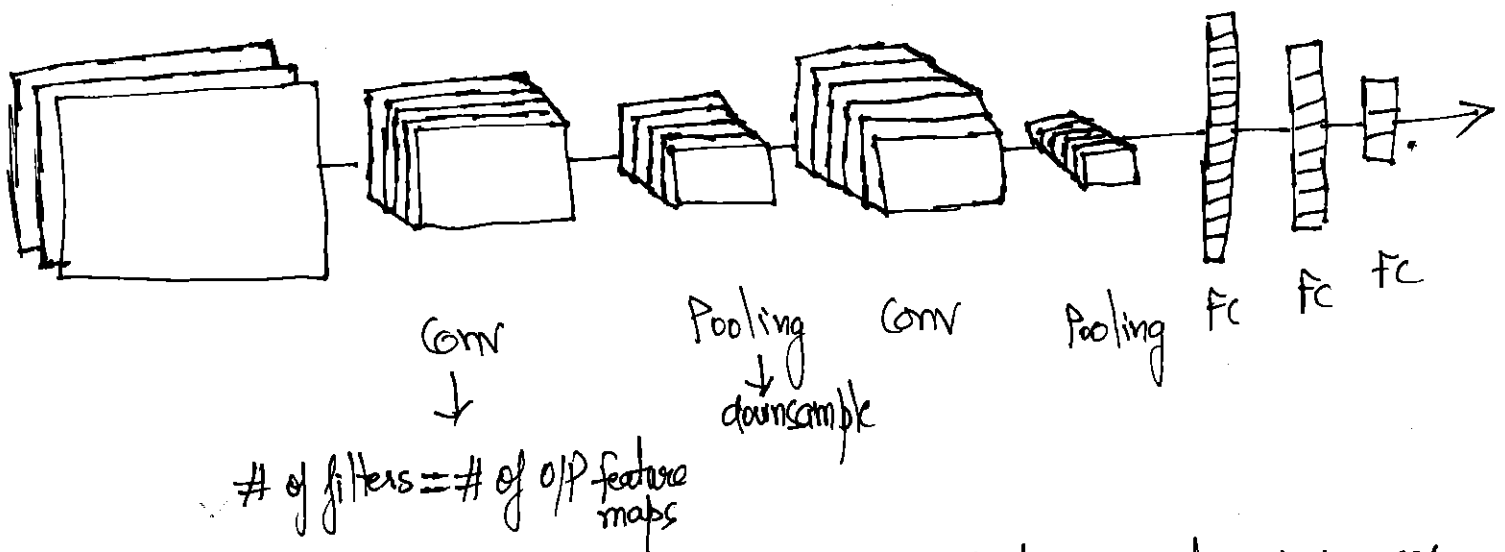
→ Pooling layer has no wts  $\Rightarrow$  all it does is aggregate the inputs using an agg. func. such as mean, max etc.

# CNN ARCHITECTURES

Typical CNN architecture involves:

- 1) (Convolutional layer<sub>(+Relu)</sub> + Pooling layer)  $\times n$  + Fully connected layer
- or
- 2) (Convolutional layer<sub>(+Relu)</sub>  $\times n$  + Pooling layer)  $\times n$  + Fully connected layer

The image gets smaller and smaller as it progresses thru the n/w but it typically gets deeper and deeper (i.e. with more feature maps)  
typically  $\downarrow$  # of filters increase as the n/w size increases



$\Rightarrow$  width and ht. decreases but # o/p feature maps increases (since # of filters typically increase when going deep)

## CONVNET CODE (KERAS)

```
model = models.Sequential()  
model.add(layers.Conv2D(32, (3,3), activation='relu',  
                        input_shape=(28,28,1)))
```

```
model.add(layers.MaxPooling2D((2,2)))
```

```
model.add(layers.Conv2D(64, (3,3), activation='relu'))
```

```
model.add(layers.MaxPooling2D((2,2)))
```

```
model.add(layers.Conv2D(64, (3,3), activation='relu'))
```

```
model.add(layers.Flatten())
```

```
model.add(layers.Dense(64, activation='relu'))
```

```
model.add(layers.Dense(10, activation='softmax'))
```

```
model.summary()
```

<u>Layer</u>	<u>O/P shape</u>	<u>Param #</u>
conv2d_1	(None, 26, 26, 32)	320 $(32 \times 3 \times 3) + \frac{32}{\text{bias}}$ # of filters    filter size    # of filters
maxpooling2d_1	(None, 13, 13, 32)	0 (No trainable params for pooling)

Conv2d-2

(None, 11, 11, 64)

18496

$(32 \times 64 \times 3 \times 3 + 64)$   
# of I/P filters from prev. conv layer   # of filters in this layer   size of filters   bias # of filters in this layer

max pooling2d-2

(None, 5, 5, 64)

0

Conv2d-3

(None, 3, 3, 64)

$64 \times 64 \times 3 \times 3 + 64$   
 $= 36928$

Flatten-1

(None, 576)  
3x3x64

0

Dense-1

(None, 64)

36928  
 $(576 \times 64 + 64)$

Dense-2

(None, 10)

650  
 $(64 \times 10 + 10)$   
I/E from prev. layer   # of nodes in this layer   bias # of nodes in this layer

Total params: 93,322

IMP: Pooling layer has zero trainable params  
= all it does is aggregate

LEARNABLE PARAMS CALC:

Inputs x outputs + biases  
↓   ↓  
i) If last layer is dense = # of nodes   # of filters x size of filters  
ii) If last layer is conv = # of filters   # of filters



→ Filters define the wt and biases

→ One feature map share wts & biases i.e. different feature maps use different parameters

→ The fact that all neurons in a feature map share the same parameters dramatically reduces the number of parameters in the model.

⇒ Once the CNN has learned to recognize a pattern in one location, it can recognize in any other location

## TYPES OF COMPUTER VISION PROBLEMS:

- i) Image Classification — Image is of cat vs dog
- ii) Object detection — In an image, what are the different objects e.g. wall, cat, vessel, shoes etc.
- iii) Semantic Segmentation — In an image, what are the boundaries of different objects e.g. self driving cars → how diff. objects are delineated

## ✓ PURPOSE OF $1 \times 1$ CONVOLUTION

Pooling layer 'typically' downsamples width & ht of I/P feature map.  
with no change in depth / # of output feature maps.

Conv. layer 'typically' increases depth with no change in width  
or ht of I/P feature map (if zero padding)

$1 \times 1$  conv. is used to downsample or reduce the depth  
or # of o/p feature maps

So,  $1 \times 1$  conv  $\rightarrow$  Channel-wise / depth-wise pooling

$\downarrow$  can be thought of as:  
dimensionality reduction  
(shrinks the number of channels)

## PURPOSE OF SKIP CONNECTIONS (RESIDUAL NETWORKS)

**SKIP CONNECTIONS:** The signal feeding into a layer is also added to the o/p of a layer located a bit higher up

The idea is if you take a "shallow" network and just stack on more layers to create a deeper network, the performance of the deeper network should be at least as good as the shallow network since the layers in b/w shallow & deep n/w should be at least able to learn the identity fn. However in practice, this was not the case  $\rightarrow$  hence residual learning or skip connections

So this motivated the use of "deep residual layers" to allow network to learn deviations from the identity layer, hence the term "residual", residual here referring to difference from identity.

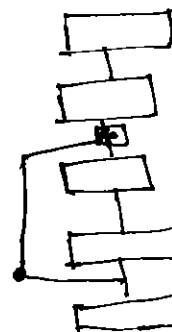
### PURPOSE:

- 1) They overcome the issue of vanishing gradients and hence help in gradient propagation



$\rightarrow$  Layer blocking back propagation

$\rightarrow$  Layer not learning



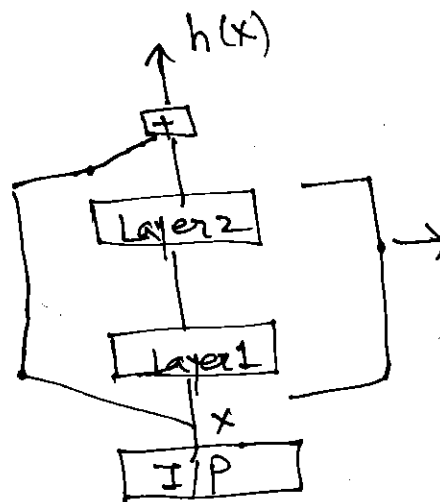
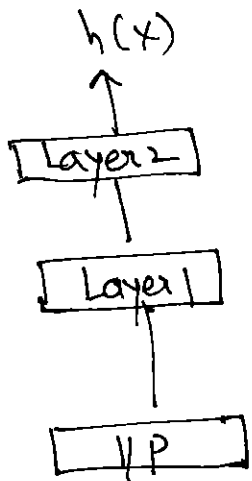
} Residual Units

$\rightarrow$  Layer starts learning

2) They make the skipped layers learn an identity function  
means the higher layer n/w will perform atleast  
as good as the lower layer n/w

## RESIDUAL LEARNING:

When training a neural n/w, the goal is to make it  
model a target function  $h(x)$ . If you add the input  $x$   
✓ to the output of n/w (i.e. you add a skip connection)  
then the n/w will be forced to model  $h(x) - x$   
rather than  $h(x)$  → residual learning



→ will be forced to  
learn  $h(x) - x$

# RNN

## Diff. b/w Feedforward Neural N/w & RNN:

Feedforward neural n/w do not maintain info about order of features in an instance  
Whereas .

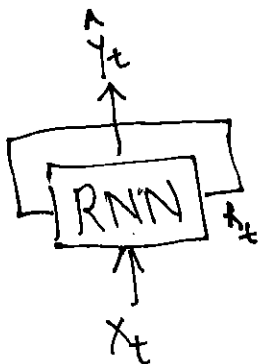
RNNs maintain info about order of elements within a sequence. RNNs assume elements within a sequence have some order but different sequences are independent.

Instance / Example in FF Neural N/w = Sequence in RNNs

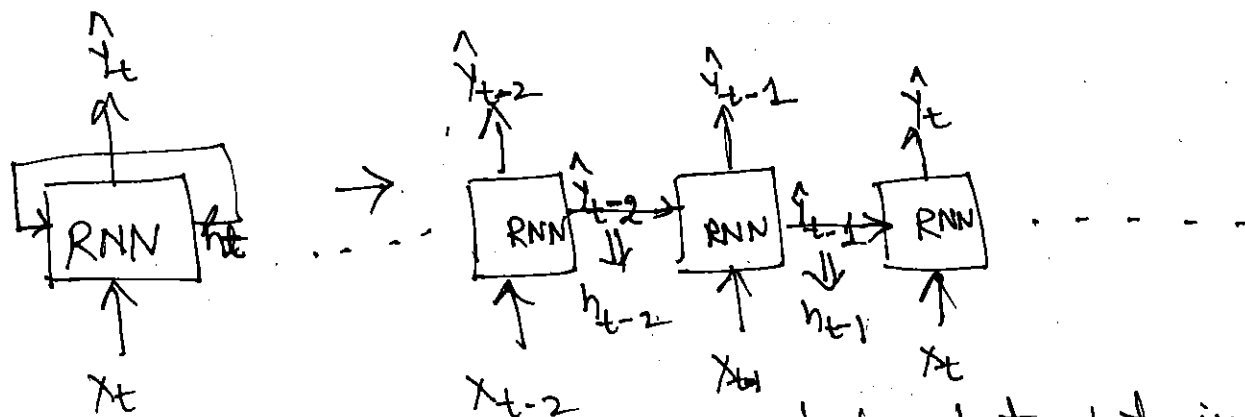
↓  
RNNs internally loop over every element in the sequence

## UNROLLING THE NETWORK THRU TIME:

Every element in a sequence is at a different time-step.  
meaning 1<sup>st</sup> element is at timestep 1  
2<sup>nd</sup> element is at timestep 2



At every timestep  $t$ , every neuron receives both the input vector  $x_t$  and output vector from previous timestep  $y_{t-1}$ .  
Since this O/P vector  $y_{t-1}$  keeps accumulating info from prev timesteps e.g.  $y_{t-2}, y_{t-3}, y_{t-4} \dots$  etc.



In basic cells, the o/p is simply equal to state but in complex cells, hidden state may be different than o/p.  
 Recurrent: Info is passed from one timesteps to next in a seq.

### RNNs:

Apply a recurrence relation at every timestep to process a seq.

$$h_t = f_w(h_{t-1}, x_t)$$

$\downarrow$  cell state       $\downarrow$  fn parameterized by  $w$        $\downarrow$  state vector       $\downarrow$  I/P at time step  $t$

RNNs maintain an internal state  $h_t$ . At every timestep they apply a function  $f$  (parameterized by a set of wts.  $w$ ) to update the state  $h$ .

This state update is dependent on prev. state  $h_{t-1}$  as well as current input.

Imp. to note: the same function (with the same set of params) are used at every timestep in a sequence

→ Just like in feedforward neural nets, wts are learnt across instances, in RNNs wts. are learnt across sequences i.e. wts and biases (params) are shared for a sequence

→ Every recurrent cell has two sets of wts: one depending of I/P  $x_t$  and other depending on prev. hidden state ( $w_h$ )

→ I/P vector :

$x_t$



update hidden state :

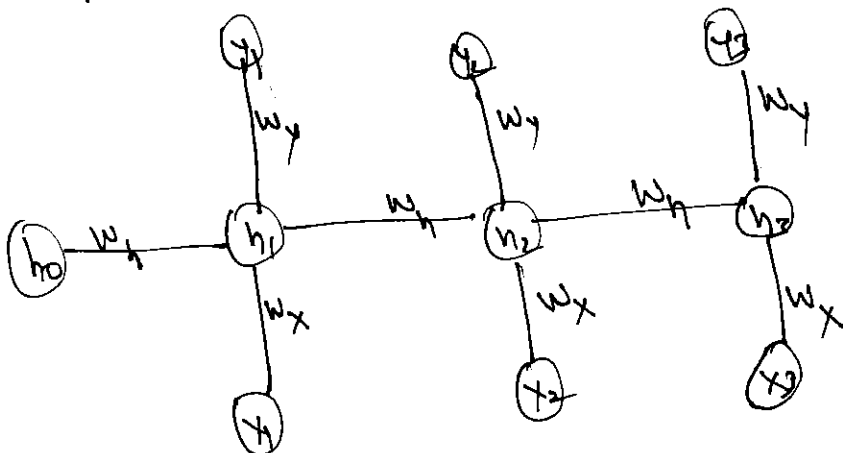
$$h_t = \tanh (w_h^T h_{t-1} + w_x^T x_t)$$



O/P vector :

$$\hat{y}_t = w_y^T h_t$$

→ We have  $w_y$  since same input can produce diff. o/p's depending on hidden state (or prev. inputs in the sequence)





# TYPES OF RNN BASED ON I/P & O/P (APPLICATIONS OF RNN)

## 1. Sequence-to-Seq (Many-to-Many)

An RNN can take a sequence of I/Ps and produce a sequence of O/Ps.

Example: Time series - you feed last N days of data and it will produce next N days predictions

## 2. Sequence-to-Vector (Many-to-One)

An RNN can take sequence of I/Ps and ignore all O/Ps except the last one.

Example: you could feed a seq. of words corresponding to a movie review and the network would o/p a sentiment score

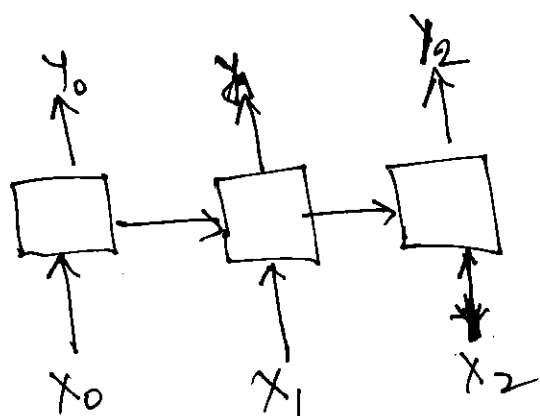
## 3. Vector-to-Sequence (One-to-Many)

You could feed the n/w the same I/P vector over and over again at each timestep and let it o/p a sequence.

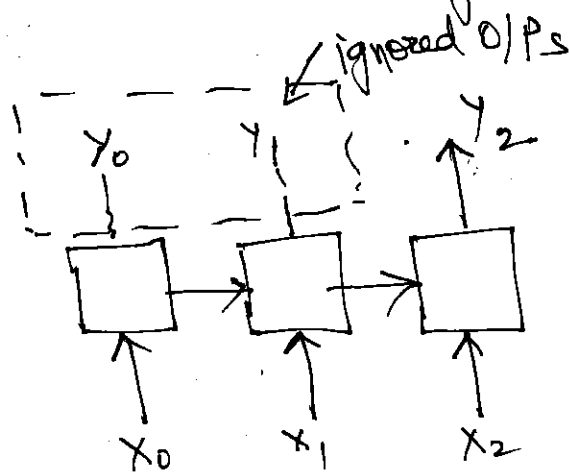
Example: I/P could be an image and O/P could be caption for that image

4. Seq-to-Vector followed by Vector-to-Seq  
↓ encoder ↓ decoder

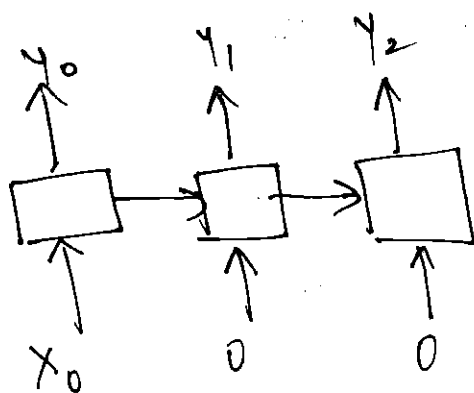
Example: Translating a sentence from one lang. to another  
You could feed the n/w a sent. in one lang. and the encoder would convert this sent. into a single vector representation, and then the decoder would decode this vector into a sentence in another lang.



Seq-to-Seq



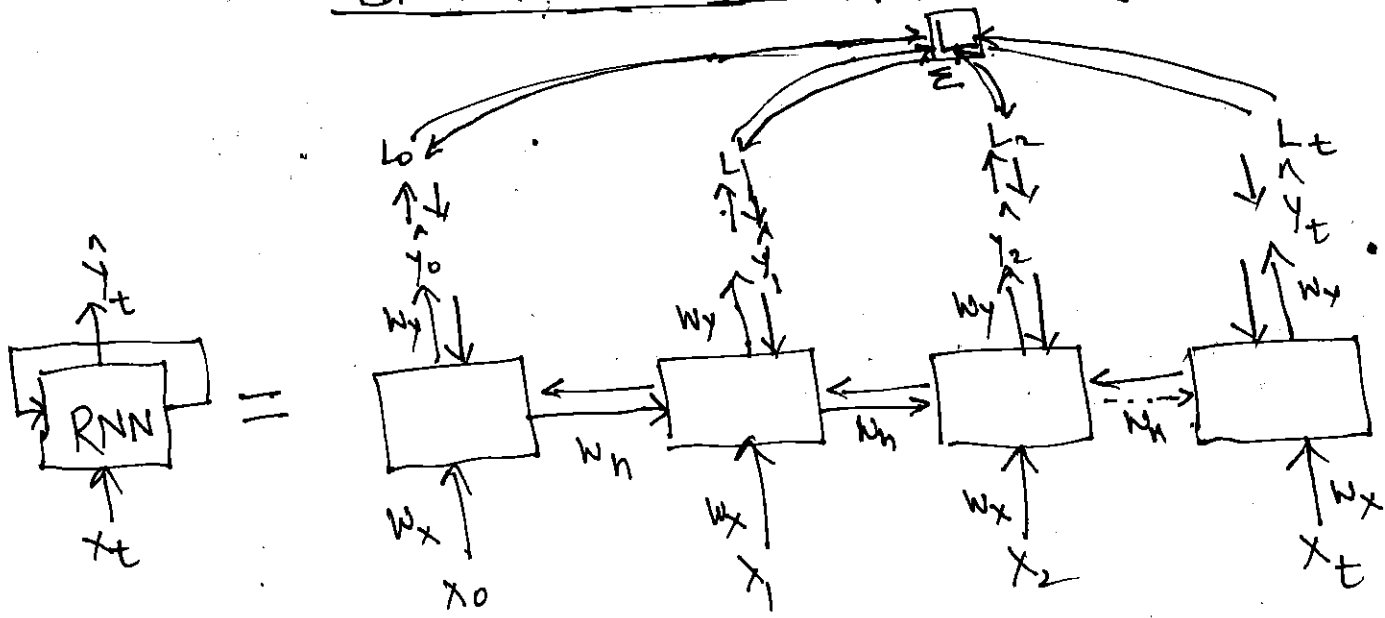
Seq-to-Vector



Vector-to-Seq

# TRAINING RNNs

## BACK-PROPAGATION THRU TIME

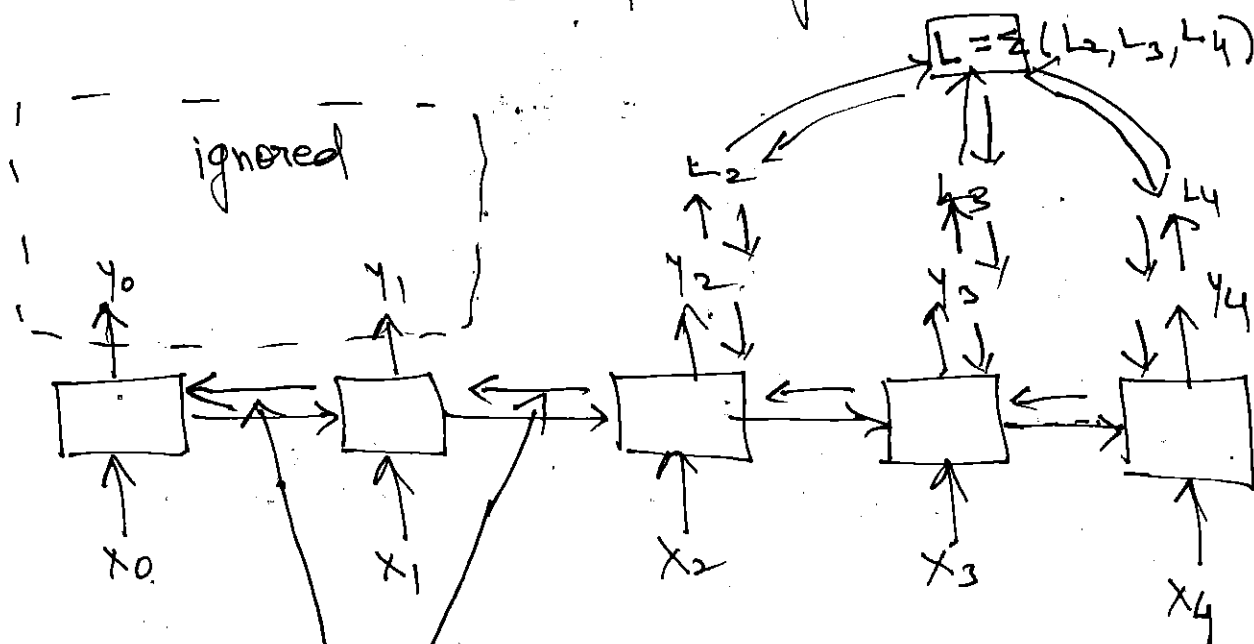


$$\text{Cost Fn} : L(\hat{y}_0, \hat{y}_1, \dots, \hat{y}_t) = \sum_{t=\text{max timestep}} (L_0, L_1, L_2, \dots, L_t)$$

This cost/loss fn may ignore some O/Ps depending on the configuration of RNN ex. if seq-to-vector conf. is used, only  $L_t$  is taken into account and rest of the O/Ps  $L_0, L_1$ , etc. are ignored. Also the error is propagated back only to these O/Ps (which the cost fn uses)

Errors are back-propagated at each timestep back upto the first timestep

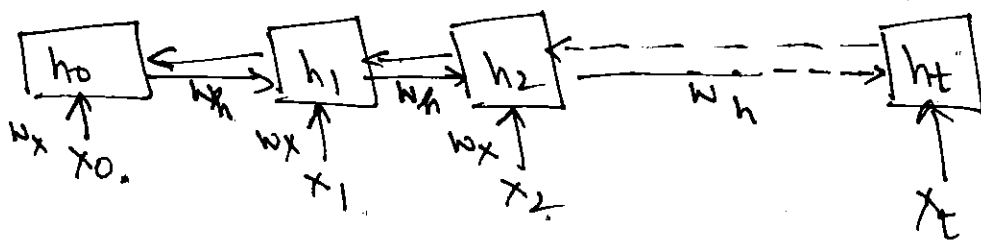
In the scenario when  $y_0$  &  $y_1$  are ignored:



These are still used even though  $y_0$  &  $y_1$  are not used to compute total loss

Since the parameters  $w_x, w_y, w_h$  are used at each time step, backpropagation will sum over all time steps

Country	1950	1960	1970	1980	1990	2000	2010	2020	2030	2040	2050
Japan	7	8	10	12	14	16	18	20	22	24	26
Germany	10	11	12	13	14	15	16	17	18	19	20
France	11	12	13	14	15	16	17	18	19	20	21
Italy	12	13	14	15	16	17	18	19	20	21	22
Spain	13	14	15	16	17	18	19	20	21	22	23
Sweden	14	15	16	17	18	19	20	21	22	23	24
United Kingdom	15	16	17	18	19	20	21	22	23	24	25
United States	16	17	18	19	20	21	22	23	24	25	26
Canada	17	18	19	20	21	22	23	24	25	26	27
South Korea	18	19	20	21	22	23	24	25	26	27	28
China	19	20	21	22	23	24	25	26	27	28	29
India	20	21	22	23	24	25	26	27	28	29	30
Indonesia	21	22	23	24	25	26	27	28	29	30	31
Philippines	22	23	24	25	26	27	28	29	30	31	32
Thailand	23	24	25	26	27	28	29	30	31	32	33
Malaysia	24	25	26	27	28	29	30	31	32	33	34
Singapore	25	26	27	28	29	30	31	32	33	34	35
South Africa	26	27	28	29	30	31	32	33	34	35	36
Argentina	27	28	29	30	31	32	33	34	35	36	37
Brazil	28	29	30	31	32	33	34	35	36	37	38
Mexico	29	30	31	32	33	34	35	36	37	38	39
Colombia	30	31	32	33	34	35	36	37	38	39	40
Venezuela	31	32	33	34	35	36	37	38	39	40	41
Chile	32	33	34	35	36	37	38	39	40	41	42
Peru	33	34	35	36	37	38	39	40	41	42	43
Ecuador	34	35	36	37	38	39	40	41	42	43	44
Bolivia	35	36	37	38	39	40	41	42	43	44	45
Paraguay	36	37	38	39	40	41	42	43	44	45	46
Uruguay	37	38	39	40	41	42	43	44	45	46	47
Puerto Rico	38	39	40	41	42	43	44	45	46	47	48
Costa Rica	39	40	41	42	43	44	45	46	47	48	49
Panama	40	41	42	43	44	45	46	47	48	49	50
Dominican Republic	41	42	43	44	45	46	47	48	49	50	51
Honduras	42	43	44	45	46	47	48	49	50	51	52
Guatemala	43	44	45	46	47	48	49	50	51	52	53
El Salvador	44	45	46	47	48	49	50	51	52	53	54
Nicaragua	45	46	47	48	49	50	51	52	53	54	55
Haiti	46	47	48	49							



Between each time step we need to perform matrix multiplication involving wt. matrix  $W_x$  &  $W_h$  (and  $W_y$  if seq to seq)

Also, cell update involves matrix multiplication with non-linear activation fn.

⇒ Computation of gradient i.e. derivative of loss w.r.t. parameters ~~through~~ all the way back to first time step requires many repeated matrix multiplication of wt. matrix + repeated use of derivative of activation fn → Problematic

✓ exploding gradient

Many values involved in this repeated multiplication  $> 1$

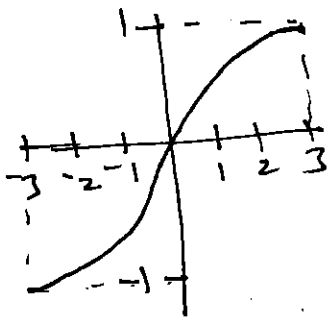
Soln: Gradient Clipping

vanishing gradient

Vanishing gradient  
Many values in this repeated multiplication  $< 1$

Soln: 1) Activation fn - <sup>relu</sup> ~~tanh~~ <sup>instead of</sup> sigmoid  
2) Wt. initialization - <sup>initialize wt to 1</sup> ~~initialize wt to 0~~  
3) Network Architecture - LSTM (long term dependence)

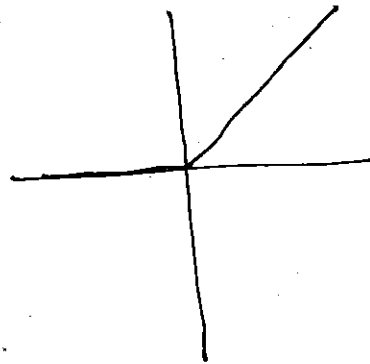
## TANH VS RELU in RNNs



TANH

Always b/w -1 & 1

↓  
Saturating Act. fn



RELU

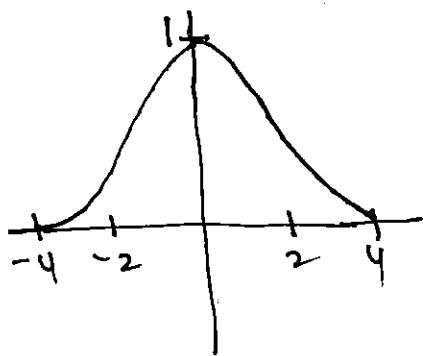
$\forall x > 0$ , unbounded

↓  
non-saturating act. fn

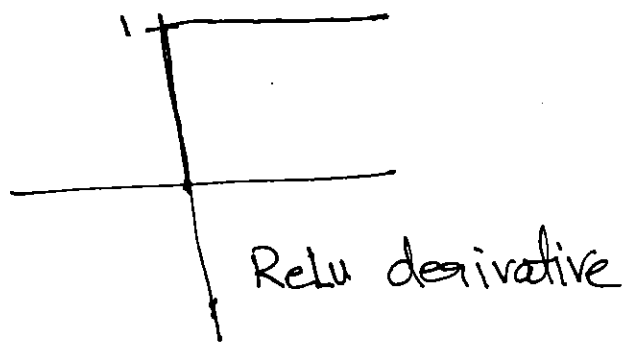
✓ To AVOID, exploding gradient problem : tanh is used as a default act. fn

Since for  $x > 0$ , relu is unbounded  
and if  $x$  goes greater than 1 then  
repeated multiplication will blow up.

However, the story is different for derivatives of  
tanh & relu and how they affect the vanishing  
gradient problem



TANH Derivative



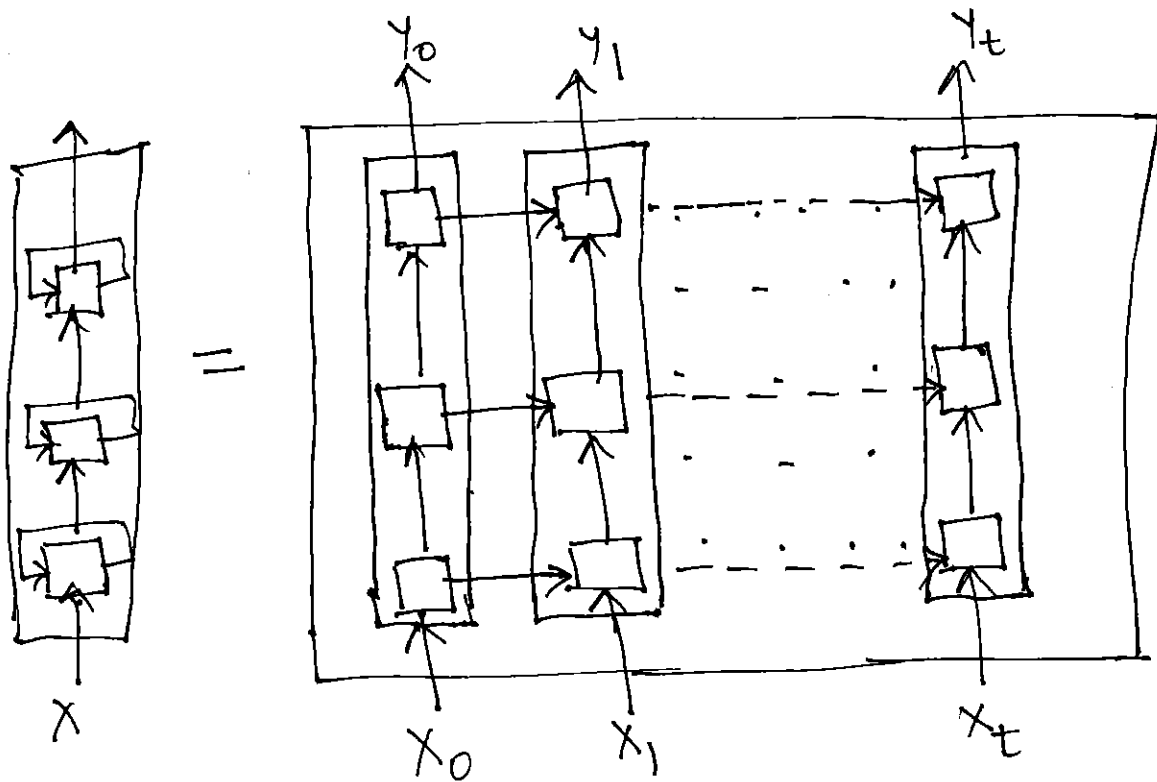
So when  $x > 0$ , relu prevents  $f'$  from shrinking the gradients

✓ To AVOID Vanishing gradient problem: relu is preferred

But since other ~~sol~~<sup>s</sup> exist for vanishing gradient problem like using LSTM etc., tanh is typically used as the activation fn.

So tanh prevents exploding gradient problem & LSTM (and other architectures e.g. GRU) prevent vanishing gradient problem.

DEEP RNN UNROLLED THRU TIME:  $\rightarrow$  stacking multiple layers



In Classical/traditional ML when modelling time series data ARIMA, MA models are used. To apply these models, time-series have to be stationary i.e. remove trend and seasonality. After the model is trained, you add trend or seasonality back to get final predictions.

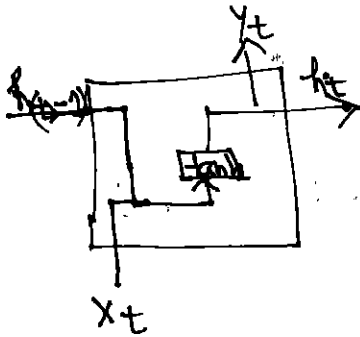
✓ When using RNN it is generally not needed to do all this. but it may improve performance in some cases.



## LSTM

- To learn long-term dependencies
- To avoid vanishing gradient problem in RNNs

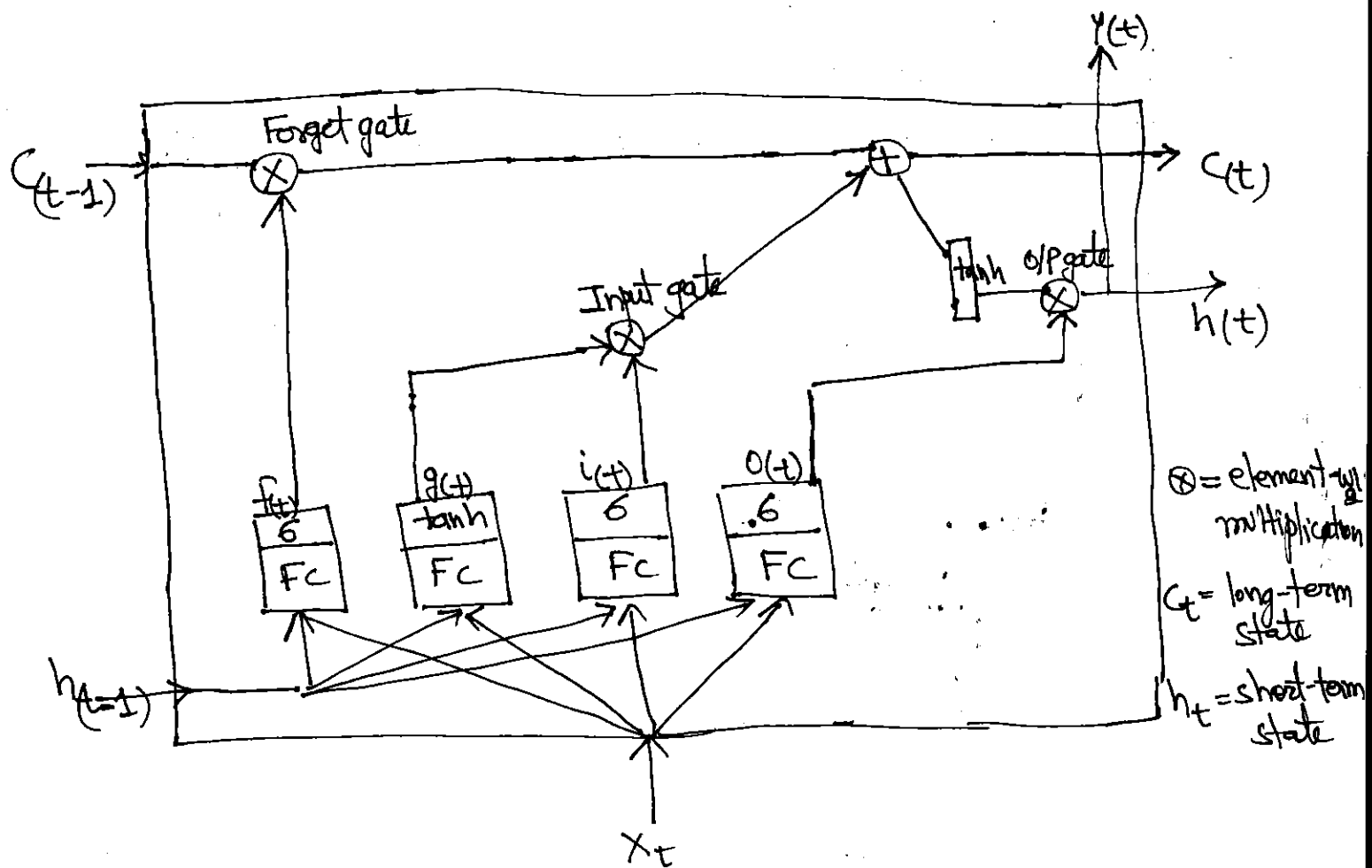
## Standard RNN



## LSTM

### ✓ Steps in LSTM:

- 1) Forget - forget / get rid of irrelevant info
- 2) Store - store relevant info from current I/P
- 3) Update - selectively update cell state (forget irrelevant info and add relevant info)
- 4) Output - output filtered version of cell state



✓ Long-term state  $C_{(t)}$  : First goes thru forget gate, dropping irr. info then adds some info that were selected by input gate  $\rightarrow C_{(t)}$

✓ Short-term state  $h_{(t)}$  :  $C_{(t)} \rightarrow \tanh \rightarrow$  filtered by o/p gate  $\rightarrow h_{(t)}$

$h_{(t)}$   
Cell's o/p at this time step

$x_t$  &  $h_{t-1}$  feeds into  $f(t)$ ,  $g(t)$ ,  $i(t)$ ,  $o(t)$

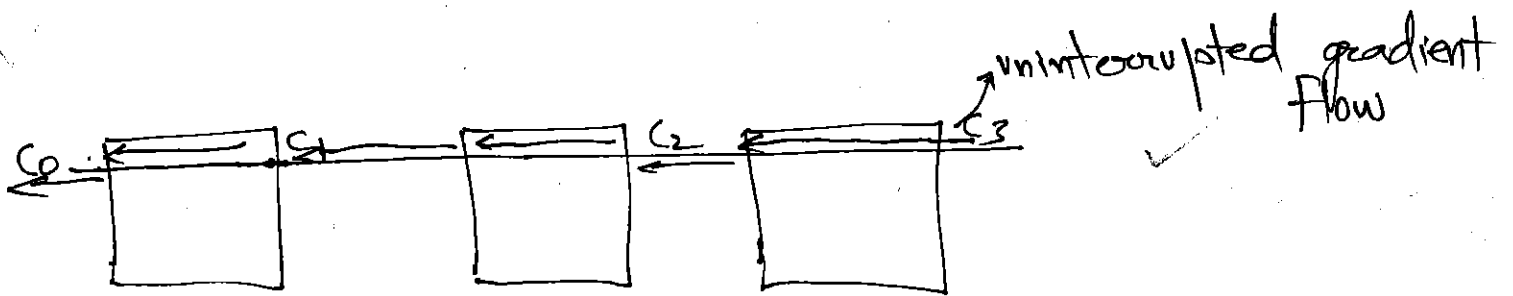
$f(t)$ ,  $i(t)$ ,  $o(t)$  are gate controllers:

### 3 GATES:

- i) forget gate is controlled by  $f(t)$ : which parts of long-term state to erase
- ii) Input gate is controlled by  $i(t)$ : which parts of  $g(t)$  should be added to long-term state
- iii) Output gate is controlled by  $o(t)$ : which parts of long-term state should be read and o/p at this time step both to  $h_t$  &  $y_t$

$g(t)$ :  $\rightarrow$  Main layer's o/p:  $g_t$  takes  $x_t$  &  $h_{t-1}$  as I/Ps and determines what part of I/P should be retained after passing thru I/P gate and added to long-term state.

\* All these gating and update mechanisms actually work to create cell state  $C$  which allows for uninterrupted flow of gradient thru time



GRU: → Similar to LSTMs but less complicated structure

→ 2 gates: reset & update (instead of 3 - input, forget, output in LSTM)

→ Reset gate: LSTM's input & forget gate functionality is done by reset gate in GRU

→ Update gate: To update the cell state

✓ GRU vs LSTM:

1. Performance-wise (detecting dependencies over long-term)  $\approx$  Similar
2. GRU is computationally efficient than LSTM  
↓  
Less complex structure of cell

## 1D CONVOLUTION LAYERS TO PROCESS SEQ

LSTMs & GRUs can process long sequences but they still have limited short-term memory and they have hard time learning long-term patterns in sequences of 100 time steps or more such as audio samples, long sentences etc.

✓ 1D convolutional layers can be used to downsample a very long seq while still retaining imp. info.  
(using strides more than 1)

1	2	3	4	5	6	7	8	9	10
---	---	---	---	---	---	---	---	---	----

 → seq (10 elements)

↓ 1D conv layer with filter size = 5  
stride = 2

1-5	3-7	5-9		
-----	-----	-----	--	--

 → seq (only 3 elements)

## EMBEDDINGS

- An embedding is a mapping from discrete objects, such as words or categorical var. value, to vectors of real numbers
- A discrete object ~~that~~ needs to be converted to some form of vector representation (e.g. one-hot encoding etc.) before embedding can be applied
- Converting a discrete object first to some vector representation ~~typ~~ results in sparse vectors (e.g. one-hot encoding) of high-dimensionality. Now when embedding space is applied to these high-dimensional sparse vectors, it results in low-dimensional dense vectors called as embeddings
- Ideally, an embedding captures some of the semantics of the input by placing semantically similar inputs close together in the embedding space. (distance <sup>b/w vectors</sup> and direction of vectors determine this semantic similarity)
- An embedding space is a low-dimensional space that translates high-dimensional vectors into low-dimensional vectors.

Discrete  
Object

us  
uk  
India  
Iran  
Taiwan  
China

one-hot  
→  
encodings

High-dimensional  
(sparse) Vectors

1  
0  
0  
0  
0  
0

0  
1  
0  
0  
0  
0

(Dense)  
Low-dimensional  
Vectors

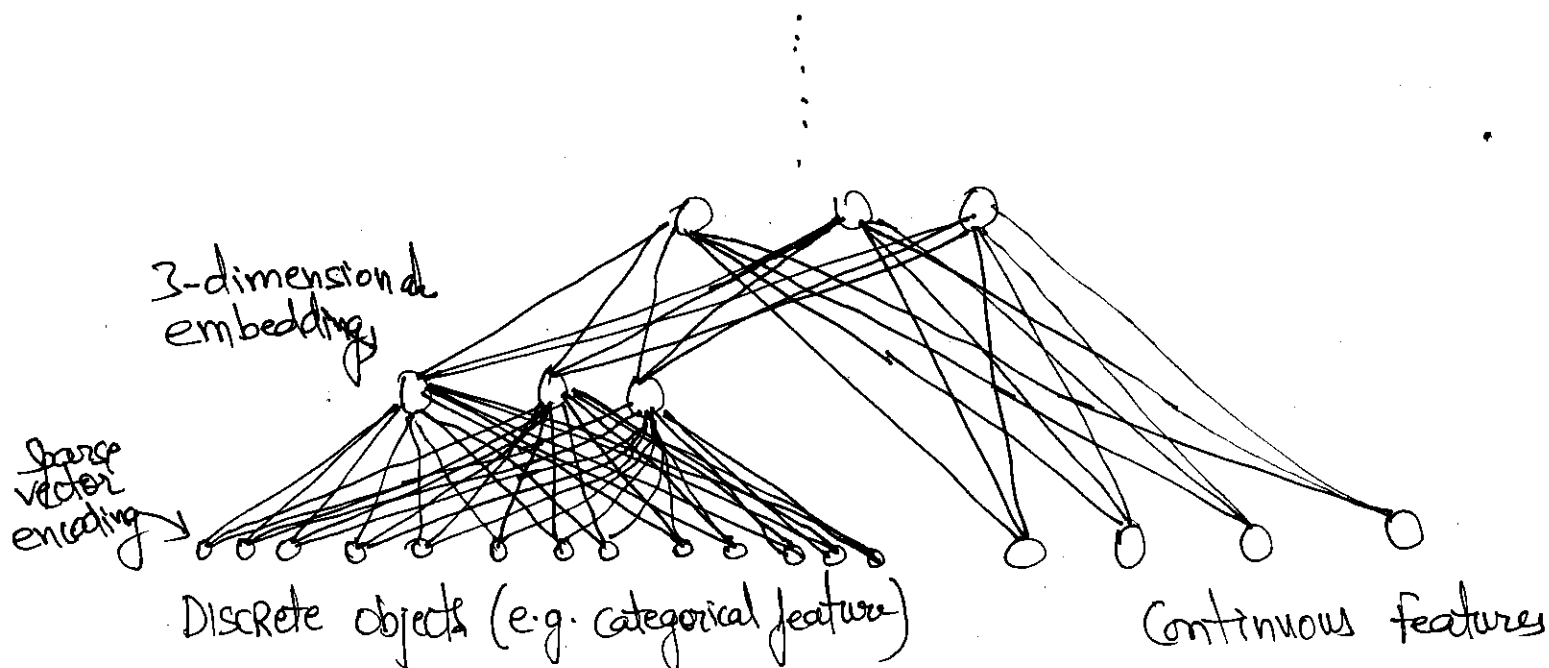
0.2  
0.6  
0.9

0.3  
0.7  
0.1

Embeddings

(Low-dimensional)  
Embedding  
space  
(Matrix)  
(Wts are learnt  
during training)

## TRAINING AN EMBEDDING AS PART OF LARGER MODEL





Embeddings can be

✓ 1. Pre-trained : e.g. the model has learnt the embedding space using some corpus (domain)

↓  
These are domain-specific  
(also called as unsupervised embedding)

Now you can take this model, ingest your discrete objects/words to this model and model outputs the embeddings.

Example: word2vec is trained on Google News dataset.

✓ 2. learnt as part of training for your task  
(supervised embedding)

Note: The concepts of word2vec can be used to learn embeddings as part of your task too.

Note: If you require any kind of distance calculation for discrete objects, embeddings is the way to go. Example, if some form of KNN is needed on categorical values  $\Rightarrow$  embeddings first then find distance

## WORD2VEC

- Combination of two tech
  - CBOW (Continuous bag of words)
  - Skip-gram

- Both of these two techniques use "shallow" neural net i.e. one hidden layer which maps an input word to a target word.

- How I/P word and o/p word is defined is determined by these two techniques

Example: Hey, this is great

CBOW: context window = 2

<u>I/P</u>	<u>O/P</u>
Hey	this
Hey	is
this	Hey
this	is
this	great
this	Hey
is	this
is	great
great	this
great	this

← Datapoint 1  
← Datapoint 2

i.e. I/P: a word (repeated)

O/P: take context window words (before & after, if they exist)

Skip-gram: (Multiple O/Ps)  
Context window = 1

✓ I/P

Hey

this

is

great

O/P

this

Hey

this

is

O/P

<padding>

is

great

<padding>

← Datapoint 1

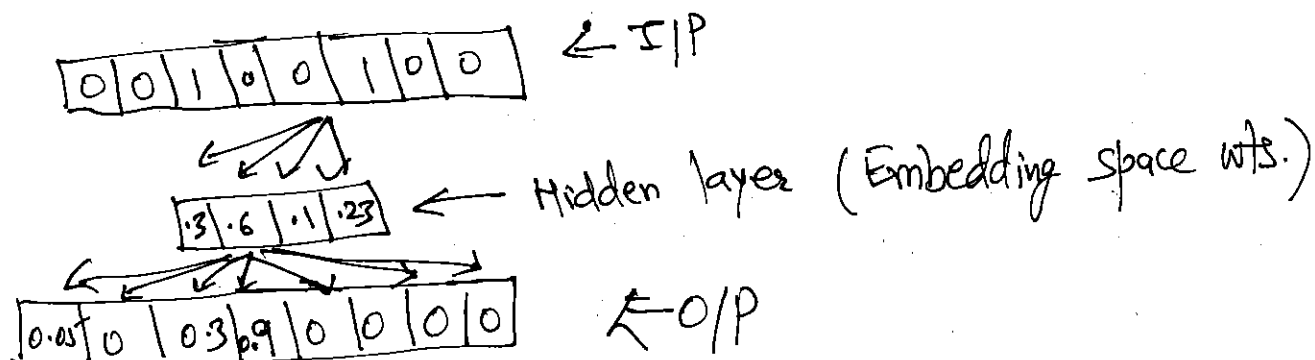
← Datapoint 2

I/P: a word

O/P: All words in context window (before & after, if they exist)

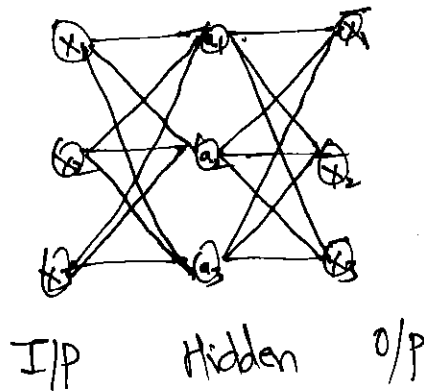
→ Since a shallow network is used, the hidden layer's wt is the embedding space

And after the model is trained, the I/P words go thru this model and the O/P of hidden layer is the embedding for the I/P word



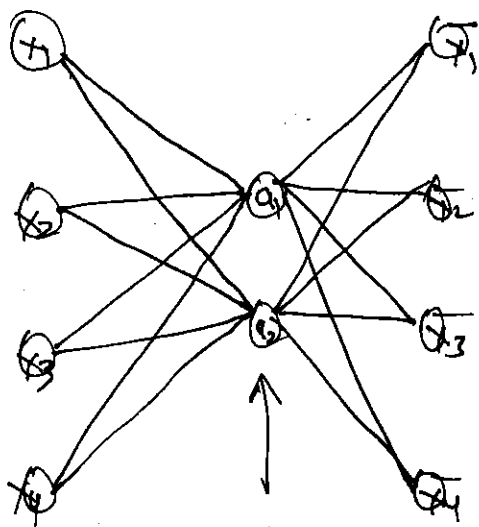
## AUTO ENCODERS (Unsupervised)

- A neural n/w architecture that imposes a bottleneck/constraint in the n/w which forces a compressed/latent representation of the original I/P.  
↓  
codings
- The bottleneck/constraint is a key attribute of our n/w design; without the presence of information bottleneck our n/w could easily learn to simply memorize the input values by passing these values along through the n/w. E.g.



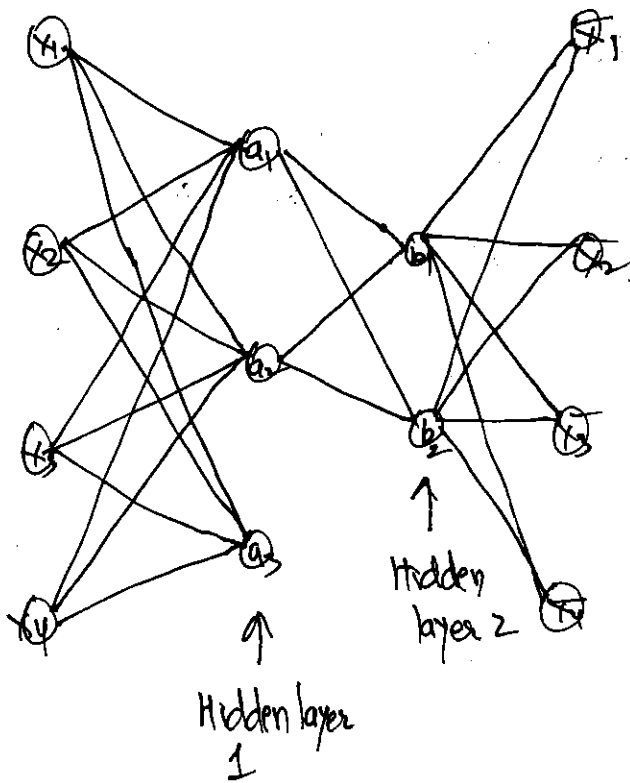
No bottleneck  
(i.e. Hidden layer size  
= I/P layer size  
= O/P layer size)

- An autoencoder looks at the input, converts them to an efficient latent representation and then spits out something that looks very close to the inputs.
- An autoencoder is always composed of two parts:
  - i) encoder: that converts I/P to latent rep
  - ii) decoder: that converts latent rep. to the O/P (= I/P)



← Undercomplete Autoencoder  
(has no explicit regularization term)

bottleneck (Hidden layer size  $\neq$  I/P or O/P layer size)  
 $\leftarrow$



← Stacked / Deep Autoencoder  
(Multiple hidden layers whose size is less than I/P or O/P layer)

↓  
Adding more hidden layers helps the autoencoder to learn more complex codings (has no explicit regularization term in theory)

Objective fn of Autoencoders:

Minimize reconstruction error + Regularizer

typically, MSE or Cross-entropy

$$L(x, \hat{x}) + \text{Regularizer}$$

when

- i) Regularizer = L1 or KL divergence → Sparse Auto-encoders (avoiding overfitting)
- ii) Regularizer = Dropout layer (or Gaussian noise) → Denoising Auto-encoders

to prevent copying of I/P to O/P

## CODE (KERAS)

```
encoder = keras.models.Sequential([keras.layers.Dense(3, input_shape=[3])])
```

```
decoder = keras.models.Sequential([keras.layers.Dense(3, input_shape=[2])])
```

```
autoencoder = keras.models.Sequential([encoder, decoder])
```

```
autoencoder.compile(loss="mse", optimizer=keras.optimizers.SGD(lr=0.1))
```

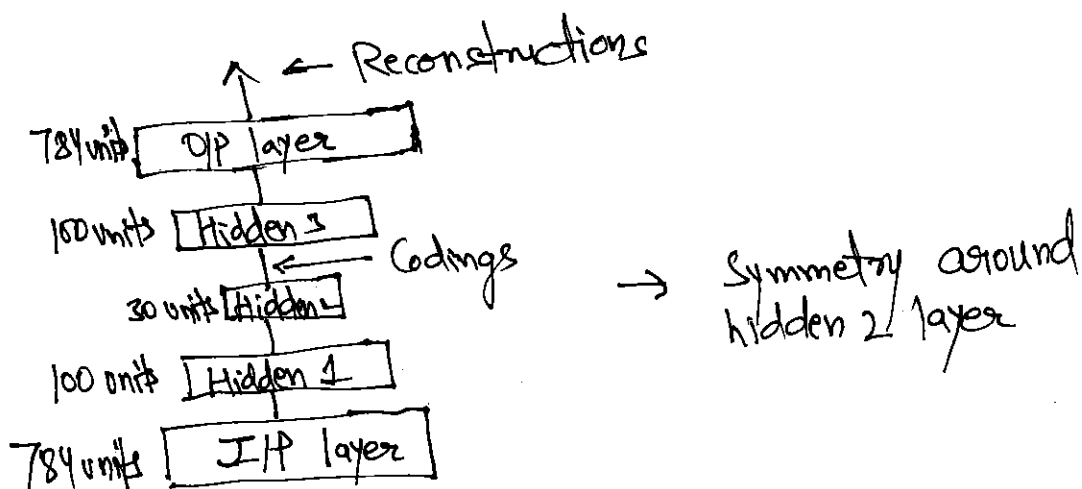
```
history = autoencoder.fit(X_train, X_train, epochs=20)
```

I/P & O/P are both X\_train

```
codings = encoder.predict(X_train)
```

↑  
latent rep<sup>s</sup> of I/P after training

## TYPICAL AUTOENCODER ARCHITECTURE:



## APPLICATIONS :

- i) Dim. Reduction / Feature Extraction  
(OIP of encoder after training)
- ii) Denoising data
- iii) Generative Models: (VAEs)

VAEs learn the parameters of prob. dist. modelling the IIP data instead of learning an arbitrary function in case of vanilla autoencoders. By sampling points from this distribution, we can use VAE as a generative model.

# PCA Vs AUTOENCODERS

Similarity:

Both attempt to discover lower dimensional rep of I/P (dimensionality reduction)

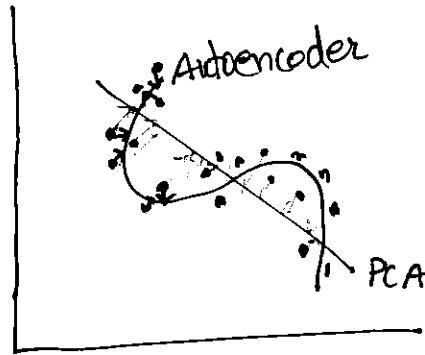
Differences:

PCA

i) Linear

Autoencoders

Non-linear



Discovers lower dimensional hyperplane which describes the original data

Capable of discovering non-linear manifolds (a manifold is a continuous non-intersecting surface)

PCA = AUTOENCODER

i) Activation in autoencoders = linear  
(recall bcz of non-linear act. fn, ~~neural~~ n/w can learn non-linear fns)

ii) Loss fn = MSE (mean squared error)



## SPARSE AUTOENCODERS

- Another type of constraint that can be used to create the bottleneck  $\rightarrow$  sparsity constraint
- In particular, we add a penalty term to the cost fn so that only a fraction of nodes become active
- This forces the autoencoder to represent each input as a comb. of smaller number of nodes and demands it to discover interesting structure in data.
- Penalty term  $\begin{cases} -L \\ -KL \text{ divergence} \end{cases}$
- This method works even if code size (size of last layer of encoder) is large since only a subset of nodes will be active at any time

## DENOISING AUTOENCODERS

Another type of constraint that can be used to force autoencoders to learn useful features  $\rightarrow$  adding random noise to its inputs and making it recover the original noise-free data (during training phase)

