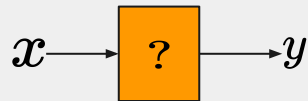




VIR - Lecture 9

Structure, recurrency, convolutions and more.

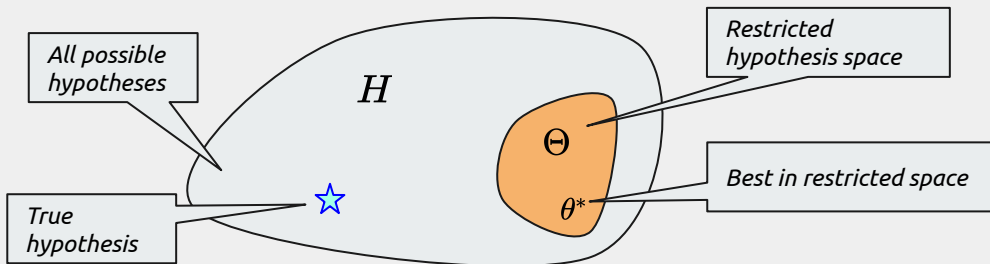
Quick learning review



Task formulation:

Given dataset $D = \{(x_1, y_1), (x_2, y_2) \dots (x_n, y_n)\}$ find hypothesis $h \in H$ which “explains” the data

Step 1: Restrict hypothesis space to something more manageable: $f_\theta : \theta \in \Theta$



Step 2: Assume that $(x_i, y_i) \sim P_{xy}$ and that it was 'reasonably' sampled (i.i.d, etc)

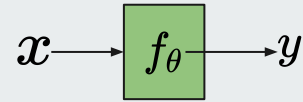
Step 3: Minimize empirical risk and hope for the best.

$$\theta^* = \arg \min_{\theta} \mathcal{L}(D, f_{\theta})$$

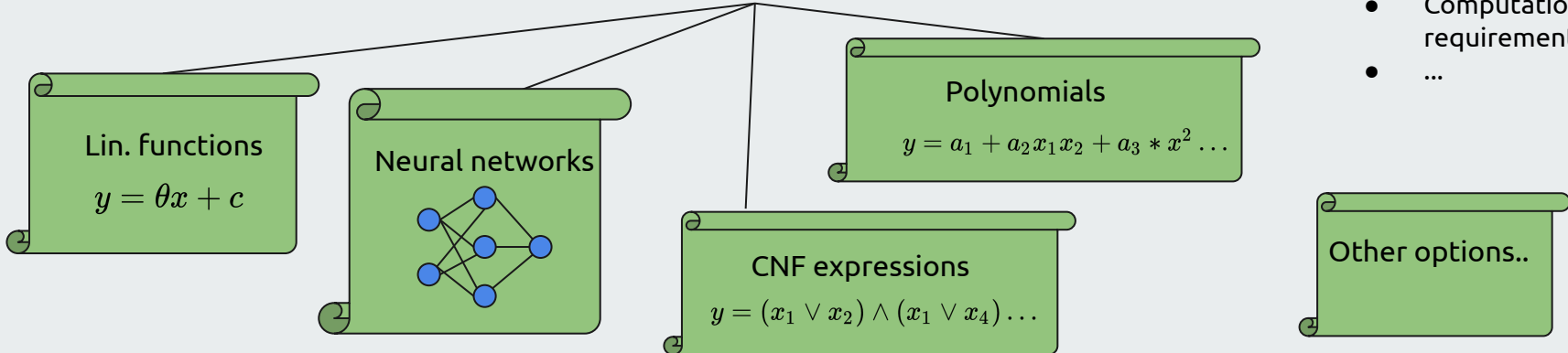
Using what you know as a 'Loss function' $\mathcal{L}(D, f_{\theta})$

Q: How do we restrict the hypothesis space?

A: Propose a suitable parameterized function.



Many potential candidates...



Criteria:

- Expressive power
- Ease of use
- Computational requirement
- ...

We will be looking mostly at neural networks and their extensions

Q: Ok, so what is the “neural network” class?

A: Loosely defined as a highly parametric functions which assume connected structure

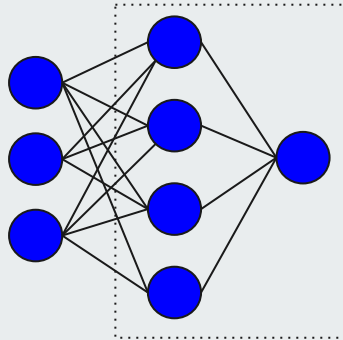
Vanilla MLP (Multi layer perceptron)

Example data:

Temp: x_1

Humidity: x_2


Pressure: x_3



y_1 Probability of rainfall tomorrow:

No assumed structure in the input data

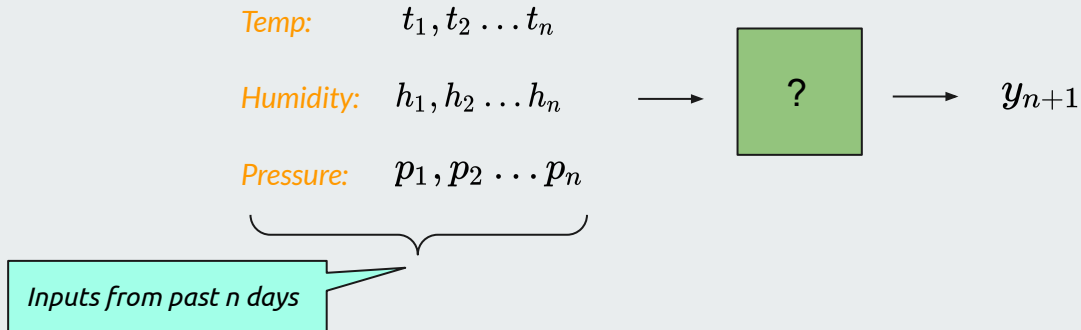
Look familiar?



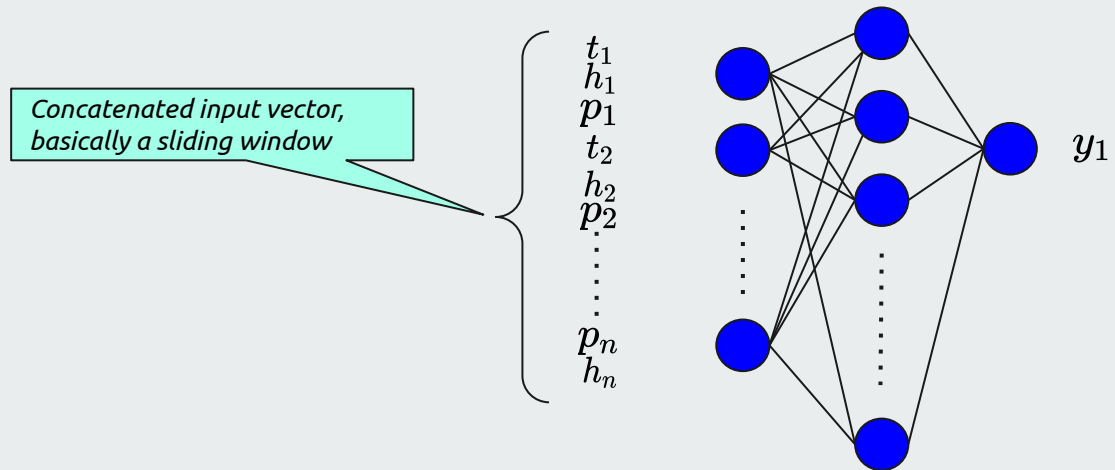
Sometimes the current state is not enough for an optimal decision. $H(Y|X)$ Is the conditional entropy

$$H(Y|X_t) < H(Y|X_t, X_{t-1}, X_{t-2} \dots) \iff \text{System not markovian}$$

So what do we do when we have a sequence of data?

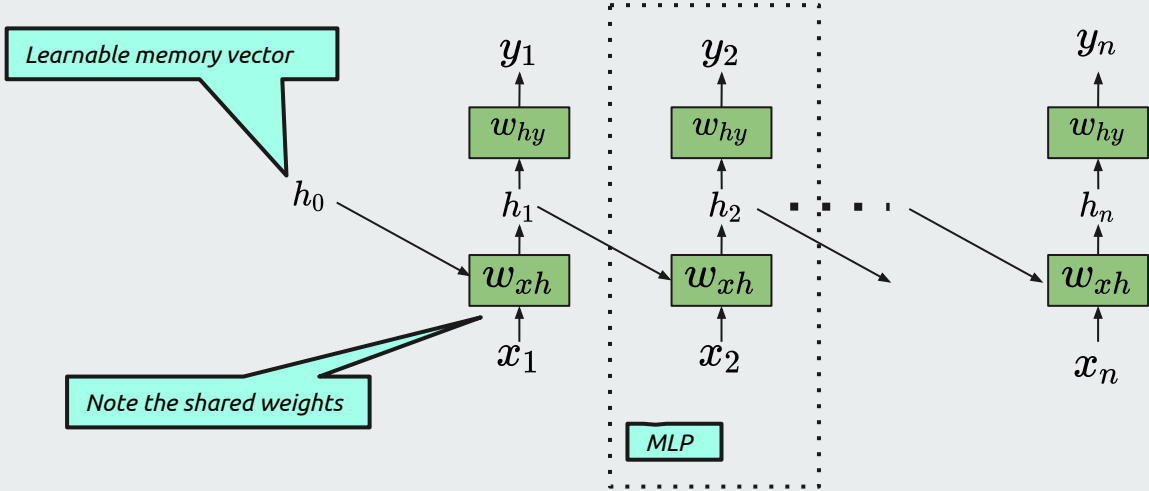


Simple solution: Concatenate everything and feed to MLP.



Will it work? Pros / cons?

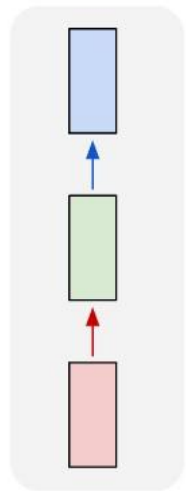
A better idea is to process data chronologically.



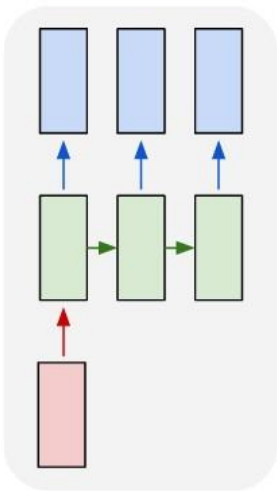
This architecture is called a recurrent neural network (RNN)

How we can make use of the RNN:

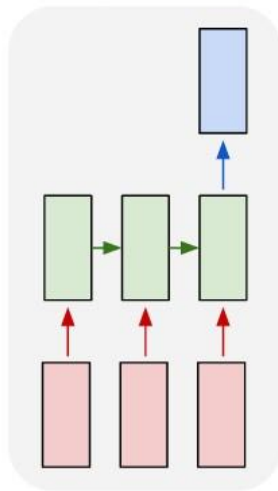
one to one



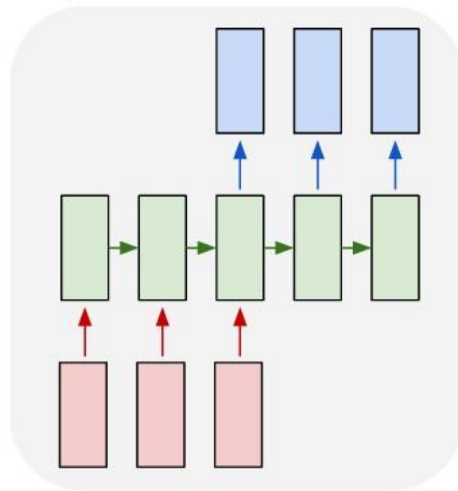
one to many



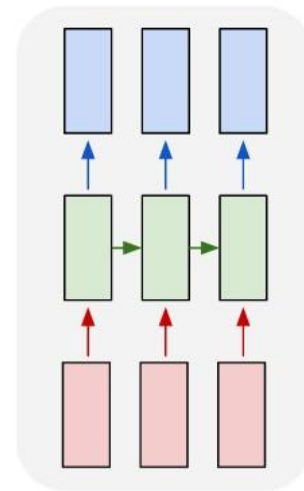
many to one



many to many



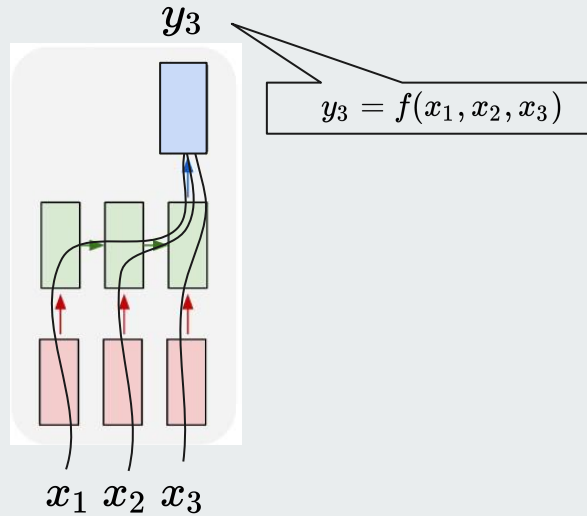
many to many



[Source: Andrej Karpathy, The Unreasonable Effectiveness of Recurrent Neural Networks]

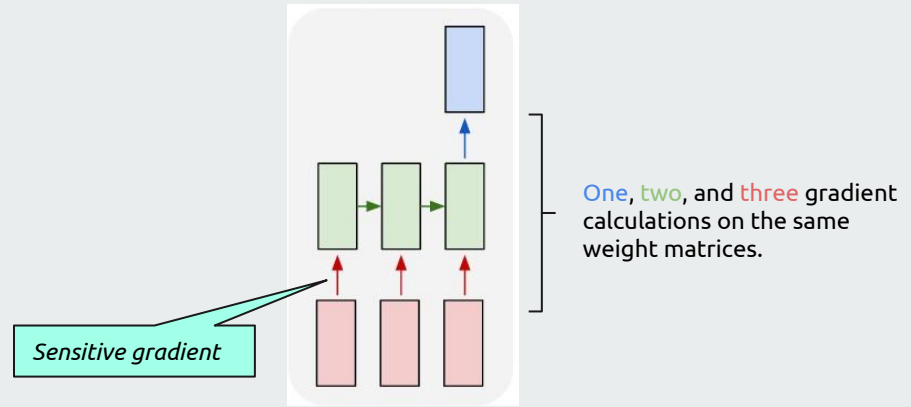
Training an RNN:

- Forward pass: Same as before, but always on the entire sequence.



$X_dim = [batchsize, seqlen, n_features]$

- Backward pass: Same as before, but with aggregated gradients



Running an RNN in Pytorch

```
import torch as T
import torch.nn as nn

x_dim, y_dim, hid_dim, num_layers = 3, 1, 8, 1

rnn = nn.RNN(input_size=x_dim, hidden_size=hid_dim,
             num_layers=num_layers, batch_first=True)

rnn_cell = nn.RNNCell(input_size=x_dim,
                      hidden_size=hid_dim)

fc = nn.Linear(hid_dim, y_dim)

batchsize = 24
seq_len = 10

# All at once
X = T.randn(batchsize, seq_len, x_dim, requires_grad=True)
H, _ = rnn(X)
Y = fc(H)

# Step by step with Cell
h = None
for i in range(seq_len):
    h = rnn_cell(X[:,i,:], h)
    y = fc(h)
```

You can specify multiple recurrent layers

Make output out of hidden layer

Works on a sequence

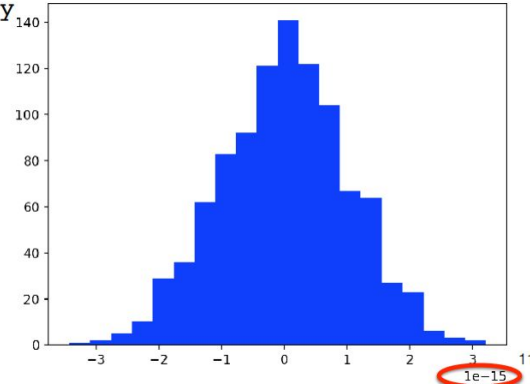
When you don't know the whole sequence in advance

The repeated multiplication problem in RNNs

RNN

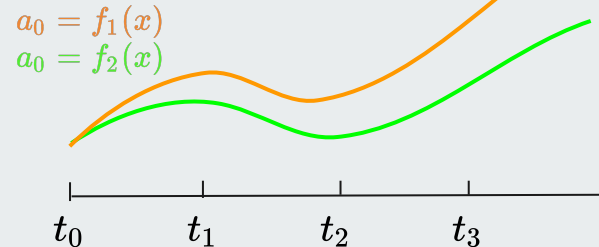
What happens to deep ~~conv~~ gradient when weights are small?

```
x = torch.randn(1000,1)
x.requires_grad_()
y=x
for i in range(30):
    weights = torch.randn(1000,1000)/100
    y = weights @ y
y.sum().backward()
x.grad
```

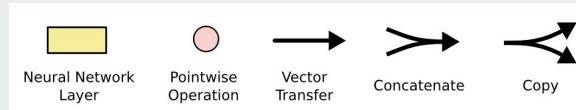
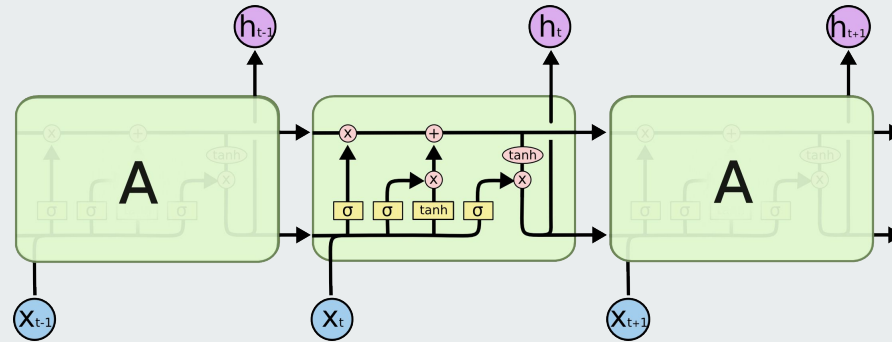


Unedited slide from VIR - Lecture IV by K. Zimmerman

- This affects both forward and backward passes
- In optimal control, this is called sensitivity to initial conditions in the 'Shooting problem'



Gated Recurrent networks: LSTM (Sepp Hochreiter, Jurgen Schmidhuber)

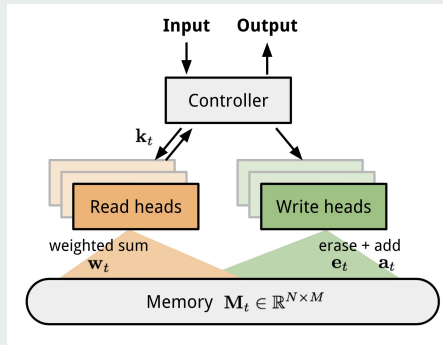


Mechanism functions takeaways

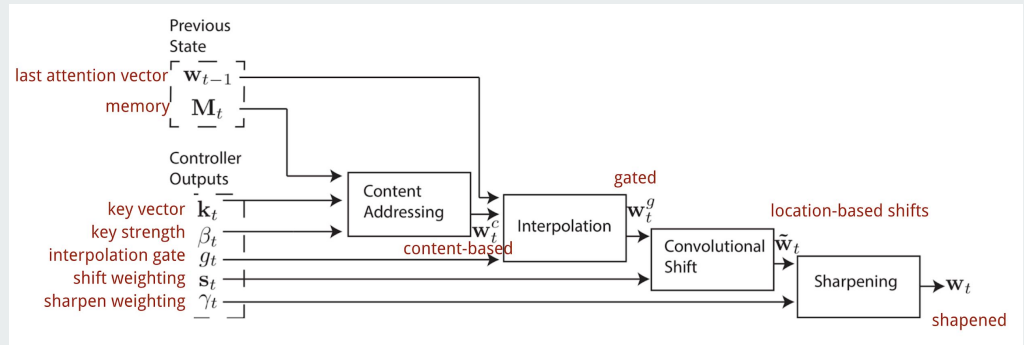
- Decide what to remember
- Decide what to forget
- Decide what to feed through

These mechanisms Ameliorate 'Forgetting' issues and gradient problems.
There are also GRU networks which are slightly simpler and do the same thing.

Recurrent network with external memory (Pretty hard to train)



Graves et al, edit: Lil-log (Lillian Weng)



Graves et al, edit: Lil-log (Lillian Weng)

Neural Turing machines (Graves et al)

What does Turing have to do with this?

Many other works in this direction:

- Neural Slam (Zhang et al.)
- Memory networks (Weston et al.)
- etc..

Examples of successes using RNNs



What color are her eyes?
What is the mustache made of?



How many slices of pizza are there?
Is this a vegetarian pizza?

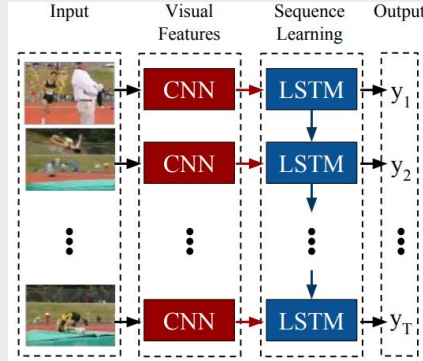


Is this person expecting company?
What is just under the tree?

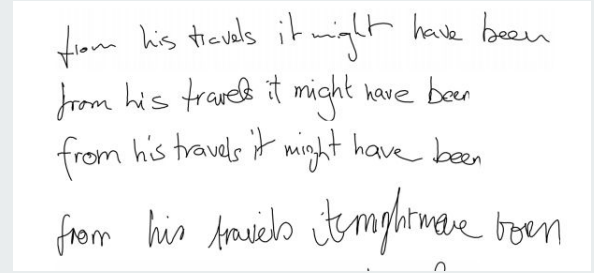


Does it appear to be rainy?
Does this person have 20/20 vision?

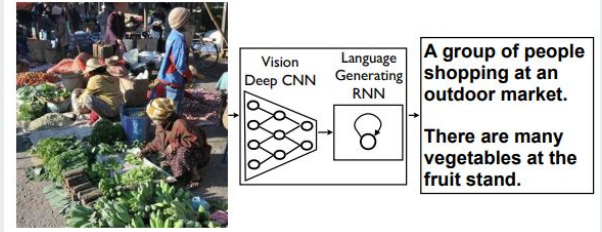
VQA: Visual Question Answering (Agrawal, Lu, Antol et al)



Recognition in video (Jeff Donahue et al.)



Generating handwriting sequences with RNN (Graves et al.)

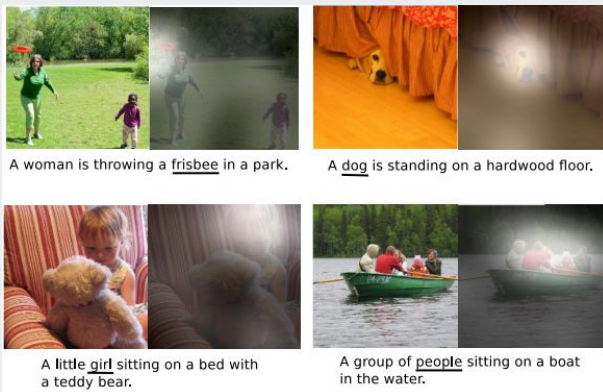


Show and Tell, Image Captioning (Vinyals et al.)

And many more...

Attention mechanisms: Focusing on relevant information

Soft attention in images:



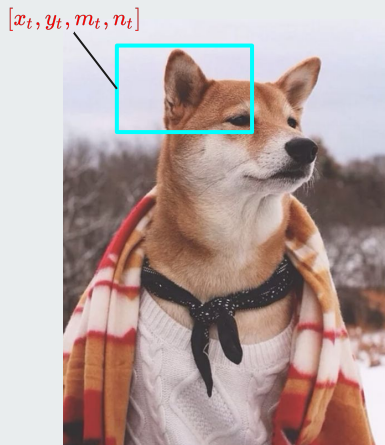
Show, Attend and Tell: Xu, et al.

Differentiable mask, convex combination of inputs:

$$y_i^{masked} = a_1 \cdot y_i \quad \text{Outputs weighted by mask coefficient}$$

$$a_1 + a_2 \dots a_n = 1$$

Hard attention in images: Retina model



Instagram [@mensweardog](#).

Soft attention in NLP:

The FBI is chasing a criminal on the run .

The **FBI** is chasing a criminal on the run .

The **FBI** **is** chasing a criminal on the run .

The **FBI** **is** **chasing** a criminal on the run .

The **FBI** **is** **chasing** a **criminal** on the run .

The **FBI** **is** **chasing** a **criminal** **on** the run .

The **FBI** **is** **chasing** a **criminal** **on** the **run** .

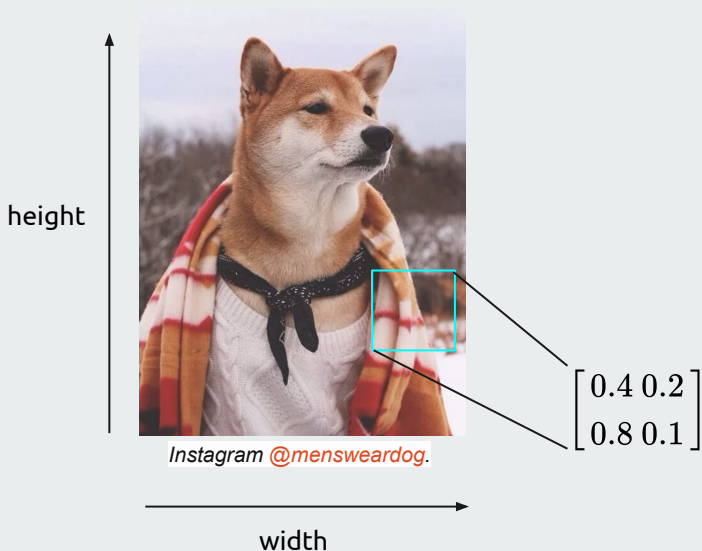
The **FBI** **is** **chasing** a **criminal** **on** the **run** .

Cheng et al., 2016



Tiny break...

Let's take a closer look at Convolution



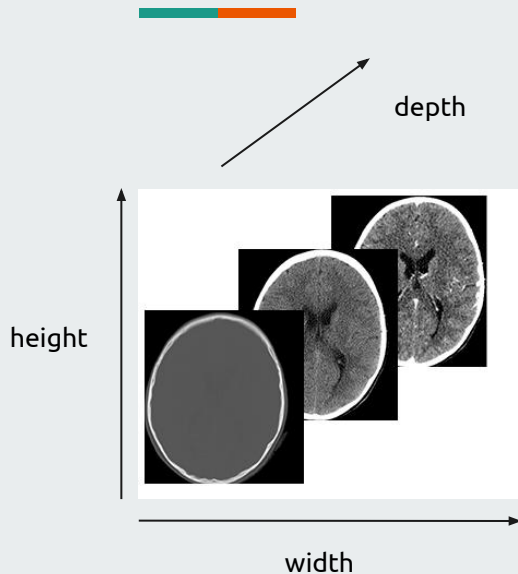
- 2D Images => 2D convolution
- Elementary filter size: $[m \times n]$
- Weight dimension: $[m \times n \times ic \times oc]$
- Filter is strided in 2 dimensions

Example of a single filter "Stamp" with filter size 2 on 1 input channel and 1 output channel:

$$\begin{matrix} c_1 & & w_1 & & b \\ \begin{bmatrix} 0.2 & 0.5 \\ 0.1 & 0.6 \end{bmatrix} & * & \begin{bmatrix} 0.1 & 0.3 \\ 0.3 & 0.4 \end{bmatrix} & + & [1] \end{matrix} = \begin{matrix} 0.2 \cdot 0.1 \\ + 0.5 \cdot 0.3 \\ + 0.1 \cdot 0.3 \\ + 0.6 \cdot 0.4 \\ + 1 = 1.44 \end{matrix}$$

If we have multiple input channels, conv separately and add

We can extend convolution to 3D



- 3D Images (sequences of images) => 3D convolution
- Elementary filter size: $[m \times n \times d]$
- Weight dimension: $[m \times n \times d \times i_c \times o_c]$
- Filter is strided in 3 dimensions

Example of a single filter “Stamp” with filter size 2 on 1 input channel and 1 output channel. Computation is same as if we are doing an image with 2 input channels

$$\begin{array}{ccc} c_1 & w_1 & b \\ \begin{bmatrix} 0.2 & 0.5 \\ 0.1 & 0.6 \end{bmatrix} & \begin{bmatrix} 0.1 & 0.3 \\ 0.3 & 0.4 \end{bmatrix} & \\ * & & + [1] = 2.14 \\ \begin{bmatrix} 0.4 & 0.4 \\ 0.7 & 0.1 \end{bmatrix} & \begin{bmatrix} 0.3 & 0.6 \\ 0.4 & 0.6 \end{bmatrix} & \end{array}$$

If we have multiple input channels, conv separately and add



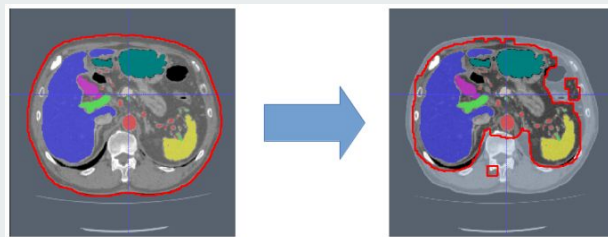
Use cases:

Video processing: Tracking, action recognition ...

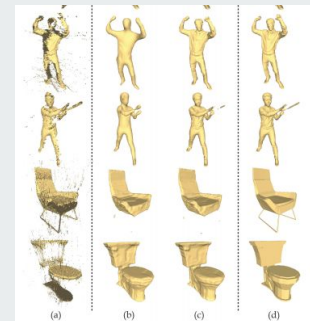


3D Convolutional Neural Networks for Human Action Recognition:
Shuiwang Ji, Wei Xu, Ming Yang

3D structure processing (3d object reconstruction, medical scans, etc.):



An application of cascaded 3D fully convolutional networks for medical image segmentation, R. Roth et al.



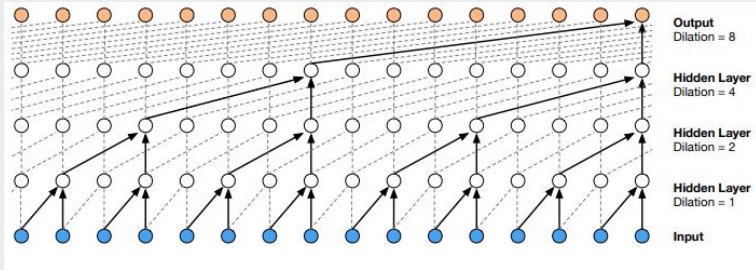
Learning to Reconstruct High-quality 3D Shapes with Cascaded Fully Convolutional Networks, Yan-Pei Cao, Zheng-Ning Liu et al.

Caveat: **What is assumed**
when using 3D conv?

We can also use convolution for 1D sequences.
(Actually a viable alternative for RNNs)



length



Wavenet. Van den Oord, et al. Deepmind.

- 1D Sequences (such as audio) => 1D convolution
- Elementary filter size: $[m \times l]$
- Weight dimension: $[m \times i_c \times o_c]$
- Filter is strided in 1 dimension

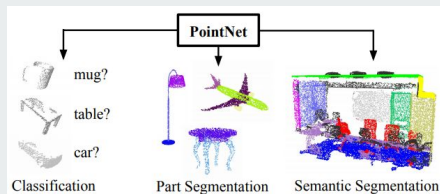
Example of a single filter “Stamp” with filter size 2 on 1 input channel and 1 output channel.

$$\begin{matrix} c_1 & & w_1 & & b \\ [0.3 & 0.1] * & [0.7 & 0.2] + & [1] = & 1.23 \end{matrix}$$

If we have multiple input channels, conv separately and add

People have been exploring how to use neural networks for other data structures

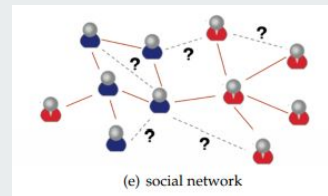
3D point clouds:



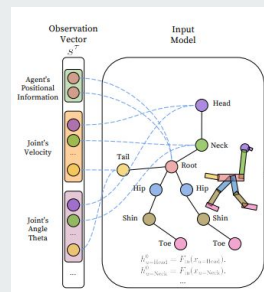
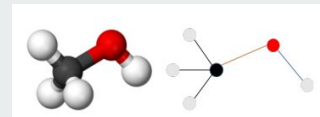
PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation. R. Qi, su et al

Graphs:

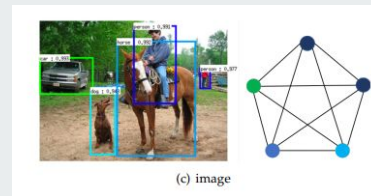
- Social connections
- Molecules
- Articulated robots
- Everything basically ...



(e) social network



Nervenet, Wang et al



(c) image

Graph Neural Networks: A Review of Methods and Applications. Zhou , Cui , Zhang et al.



Fin.
Have fun with your
semestral projects.