

11-785 CNNs review

$$t = 1 : T-K+1$$

$x(:, t:t+K-1) = x_{\text{seg}}$ i.e. a matrix? - what is dimensionality? (7)

$$y(t) = \text{MLP}(x_{\text{segment}})$$

Recall that previously

we have (x_i, y_i) $y_i \in \mathbb{R}^L$ or $\{0, 1\}$

$$\boxed{x_i \in \mathbb{R}^D}$$

- now, because sequential, each

- we have

- bit stamped by this.

- this is the 1D case?

- refresh on 1D:-

we have n input-output pairs:-

$$(x_i, y_i) \text{ i.e.}$$

- are we now saying we have
N number of $(D \times T)$ sequences? in the 1D case?

$$x_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{iT} \end{bmatrix}$$

$$\underbrace{t, t+1, \dots, t+K}_{(T)} \quad 0, 1$$

- can we scan each?

$$\text{so if we do } x^{(i)} \in \mathbb{R}^{D \times T} \quad i=1, \dots, N, \text{ we have}$$

$$\cancel{x_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{iT} \end{bmatrix}} \quad x^{(i)} = \begin{bmatrix} x_{i1}^{(i)} & \dots & x_{11}^{(i)} \\ x_{i2}^{(i)} & \dots & x_{12}^{(i)} \\ \vdots & \ddots & \vdots \\ x_{iT}^{(i)} & \dots & x_{1T}^{(i)} \end{bmatrix} = \begin{bmatrix} | & | & | & | \\ x_{:,1}^{(i)} & x_{:,2}^{(i)} & \dots & x_{:,T}^{(i)} \\ | & | & \ddots & | \end{bmatrix}$$

$$\text{then} \quad x_{\text{segment}} = x(:, t:t+K-1) \quad \text{i.e. } x^{(i)} = \begin{bmatrix} | & | & | \\ x_{:,1}^{(i)} & x_{:,2}^{(i)} & \dots & x_{:,t+K-1}^{(i)} \\ | & | & \ddots & | \end{bmatrix}$$

$$y(t) = \text{MLP}(x_{\text{segment}})$$

$$\text{so } x_{\text{segment}} \in \mathbb{R}^{D \times K}$$

$t=t+K-1$
i.e. K vectors

so what we do is

$$\underline{X}^{(i)} = \begin{bmatrix} 1 & | & \\ x_{i,1} & \dots & x_{i,T} \\ 1 & | & \end{bmatrix} \quad \underline{X}^{(i)} \in \mathbb{R}^{3 \times T}$$

0,

$$t+K-1 = t(K-1)$$

tt

If $\underline{X}[:, t:t+K-1] \text{? } \checkmark$ if $K=2$ and $t=0$
then $0:1 \rightarrow 1$ mit if

$$\underbrace{x_t, x_{t+1}, x_{t+2}, x_{t+3}}_4$$

- you are slicing using a window!

- so $y(t) =$

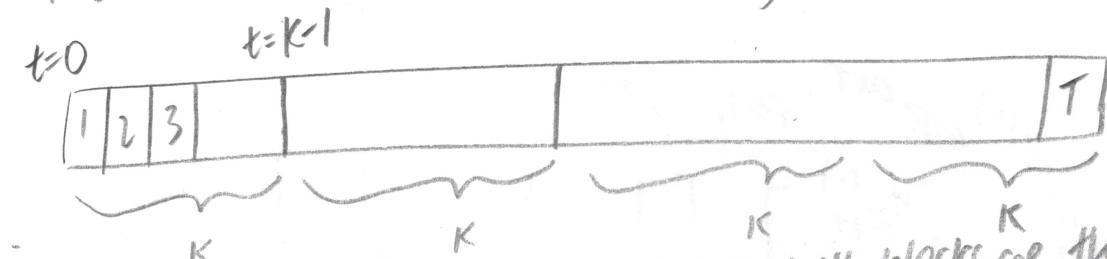
$$T-K+1 : (T-K+1) + K - 1$$

$$T-K+1 : T$$

- dimensionality of softmax output? each $y(t) \in \mathbb{R}^L$ where
 L is no. of layers in units
in final layer (or no. of classes)

- now about a giant MLP?

$y = \text{giantMLP}(\underline{X})$ - trying to clarify



- for a length T sequence, how many K -width blocks are there?

- there are

- we can rewrite multiplicatively as $T = K \cdot m$ where $m = \frac{T}{K} \in \mathbb{Z}$

- we pass each of the K -width blocks into y .

- so we will have

$$- \text{this is now sorted } N = \left\lfloor \frac{T-K}{S} \right\rfloor + 1.$$

- still not clear what
the dimensionality of y is.

Hence $y \in \mathbb{R}^{N \times (T-K+1)}$?

- If we have a K -class classification problem, then each $y_i \in \mathbb{R}^M$,
- what do we want as output?
- consider 2-class classification first.
- Prior to Yes/No, we have softmax output over 2 classes
- Recall we are going through each sliding window, of which there are $T-K+1$ sliding windows and declaring whether that window contains a 'helocar'.
- Hence $Y \in \mathbb{R}^{M \times (T-K+1)}$.

-
- 10/06/21
 - How do we have softmax from this?
 - or a giant MLP?
 - $Y \in \mathbb{R}^{M \times (T-K+1)}$
 - $\hat{y} = \text{giantMLP}(X)$ where $X = [X[:, 1:K], \dots, X[:, T-K+1:T]]$
 - and $X \in \mathbb{R}^{D \times T}$, $X[:, 1:K] \in \mathbb{R}^{D \times K}$
 - Generalisation to 2D case after 1D case is simple.
 - Need to figure out exactly what is going on with the block matrix parameter sharing.
 - What is kernel width?
 - I think I'm going to watch all the lectures again. on CNN
 - Shared parameter network

L8-CNNs: supplementary review

for the assignment HW2pt1, need to further clarify:-

1) link between scanning MLP and CNN

2) shared parameter network

3) counting parameters

4) CNN terminology

- stride
- channel
- kernel width
- \ filters

shared parameter networks

see S.2019 slide 58 lectures.

The notion of a shared parameter network is that all subnets in slide 58 of S.2019 are identical.

That is the architecture (weights and biases dimensions, activations, final out) is the same.

And also the values of the parameters in each subnet should also be ident.

And also on the fact that one update of one copy of the subnet (i.e. its parameters using gradient descent after backprop) must equally update all copies

every blue-circled subnet in slide 58 is identical

The following shows shared parameters for a simple network (not sure if scanning)

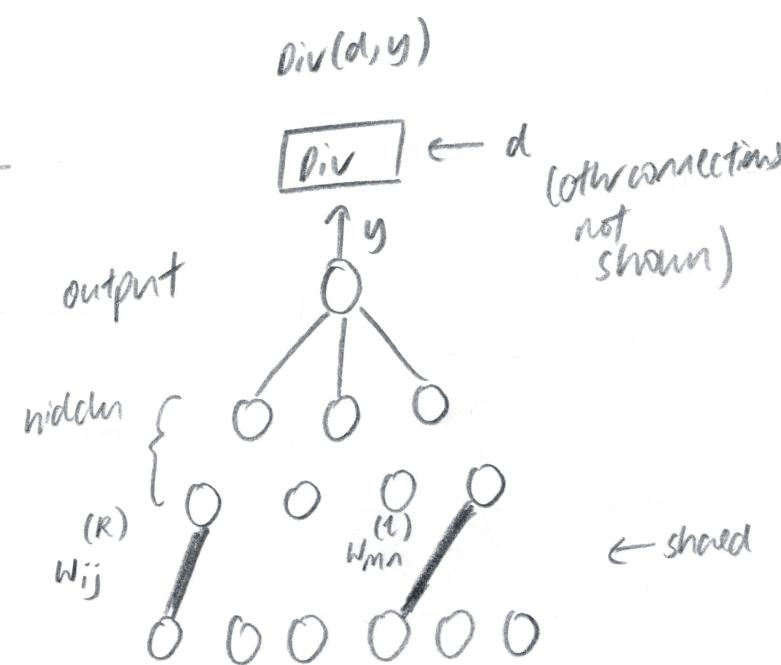
learning in shared param networks (S59)

shared weights

for example, if

$$w_{ij}^{(R)} = w_{mn}^{(1)} = w^s$$

that they are shared means that they are identical, and can be referred to using common notation w^s



It is then stated
that:-

for any training instance \mathbf{x} , a small perturbation of w^s 's perturbs $w_{ij}^{(k)}$ and $w_{mn}^{(l)}$ identically.

Each of these perturbations will influence divergence $\text{Div}(\mathbf{d}, \mathbf{y})$.

WL: I think 'perturbation' here refers to the act of updating the shared weights w^s using gradient descent and backprop.

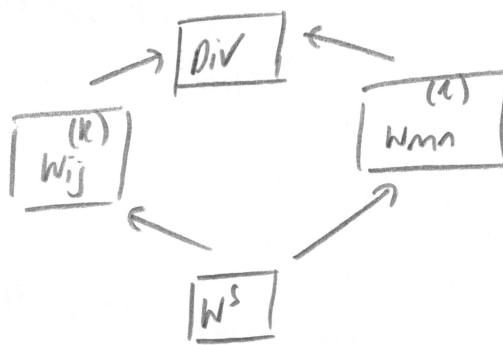
computing the divergence
of shared params (S59)

influence diagram

consider computing
derivative of loss/divergence
with respect to
shared param w^s .

chain rule:
$$\frac{d\text{Div}}{dw^s} = \frac{\partial \text{Div}}{\partial w_{ij}^{(k)}} \cdot \frac{\partial w_{ij}^{(k)}}{\partial w^s} + \frac{\partial \text{Div}}{\partial w_{mn}^{(l)}} \cdot \frac{\partial w_{mn}^{(l)}}{\partial w^s}$$
 because $\frac{\partial w_{ij}^{(k)}}{\partial w^s} = \frac{\partial w_{ij}^{(k)}}{\partial w_{ij}^{(k)}} = 1$.

$$= \frac{\partial \text{Div}}{\partial w_{ij}^{(k)}} + \frac{\partial \text{Div}}{\partial w_{mn}^{(l)}}$$



computing the
divergence of
shared parameters (S63)

key detail here :-

S - the set of edges - this is some kind of index set

$$S = \{e_1, e_2, \dots, e_N\}$$

w^s is the common weight i.e. shared parameter of the set

In S63, note that

$$\frac{d\text{Div}}{dw^s} = \sum_{e \in S} \frac{d\text{Div}}{dw^e}$$

extends previous formula.

The key idea:

Training networks with shared params (SB4)

Note that BR's notation

for layer K weight matrix is W_K .

All we do here is update every shared parameter weight.

Preliminary pseudocode

gradient descent:

- initialise all weights W_1, W_2, \dots, W_K

- DO:

- for every set S :

- compute $\nabla_{WS} \text{LOSS} = \frac{d\text{LOSS}}{dW^S}$

- compute $W^S \leftarrow W^S - \eta \nabla_{WS} \text{LOSS}^T$

- for every $(k, i, j) \in S$, update:-

- $w_{ij}^{(k)} = w^S$

- until loss has converged.

(i) The idea is that we iterate through all sets of shared indices e.g. subnets. S_1, S_2, \dots (however many shared parameter sets there are).

- update the shared parameters for that particular set S i.e. W^S

(ii) Can view this as updating weights in a single network, but updating each weight in shared parameter network amounts to updating multiple copies of each weight identically

(iii) $\nabla_{WS} \text{LOSS} = \frac{d\text{LOSS}}{dW^S} \xrightarrow{\text{by}} \text{For every training instance } X$

for every set S :-

for every $(k, i, j) \in S$:-

$$\nabla_{WS} \text{DIV} \doteq \frac{\partial \text{DIV}}{\partial w_{ij}^{(k)}}$$

computed via backprop.

$$\nabla_{WS} \text{LOSS} \doteq \nabla_{WS} \text{DIV}.$$

- That is, use previous formula $\frac{\partial \text{Div}}{\partial w^S} = \sum_{e \in S} \frac{\partial \text{Div}}{\partial w^e}$ as a subroutine in main gradient descent pseudocode.
- Looking at slide 63:
- we compute using backprop: $\frac{\partial \text{Div}}{\partial w^e} = \frac{\partial \text{Div}}{\partial w_{i,j}^{(k)}} \text{ for every } e \in S$ (i.e. compute derivative of divergence w.r.t. red edges)
- we then sum these to get the derivative of divergence w.r.t. shared weights :- $\nabla_{w^S} \text{Div} \doteq \frac{\partial \text{Div}}{\partial w^e} \text{ or } \nabla_{w^S} \text{Div} \doteq \frac{\partial \text{Div}}{\partial w_{i,j}^{(k)}}$
- However, question that remains is (possibly notation) is how or why $\nabla_{w^S} \text{Loss} \doteq \nabla_{w^S} \text{Div}$ or in S.2019 slides $\nabla_{w^S} \text{Err} \doteq \nabla_{w^S} \text{Div}$
- I think BR is referring to $\nabla_{w^S} \text{Div}$ as per-training-instance gradient, whilst $\nabla_{w^S} \text{Loss}$ is per-batch or per-mini-batch gradient (i.e. additive sum).
- Yes your understanding is completely correct
- NB: The precise nature of computations here is still not entirely clear
- ② Whilst it seems that we are saying updating w^S means we update all w^e where $e \in S$ i.e. update all shared weights w^S i.e. all multiple copies of the same weights in some set S
- ③ However does that in general mean that $\frac{\partial \text{Div}}{\partial w^e}$ in $\frac{\partial \text{Div}}{\partial w^S} = \sum_{e \in S} \frac{\partial \text{Div}}{\partial w^e}$ will be the same?
- My guess is no, because we are feeding/using a different window of 1D sequence X at each sliding window.

Key idea:

- (6) : Modify backpropagation rules to combine gradients from parameters that share the same values.
- scanning order \rightarrow 1. pass window by window through entire network.
2. pass window by window through neuron by neuron
(see orig. notes or review)
- can view it as column-wise or row-wise
- (6) - MLP scanning uses the latter order of operation
i.e. for $l=1:L$ (vectorized). - see slide 116 scanning order
for $t=1:T+l-1$

- That is run through layers one-by-one; and each layer scans
- This is precursor to scanning in CNNs
- Idea is you can switch around for loops without affecting code; even though it alters intuition of what is going on.

- I'm a bit rusty on the hierarchical explanation of scanning
- That is, from slide 117 onwards i.e. "Scanning in 2D: a closer look".
- Use the pseudocode as the foundation; overlay the intuition on top of that
- That is, look at 2D scanning pseudocode, and the re-ordered computations; overlay intuition on top of that.
- And also from slide 162 onwards, "Distributing the scan."
- See in particular, see slide 193 "Reordering the computation."
- It seems the convolutional part comes from the fact that maps scanned by each layer are a different size

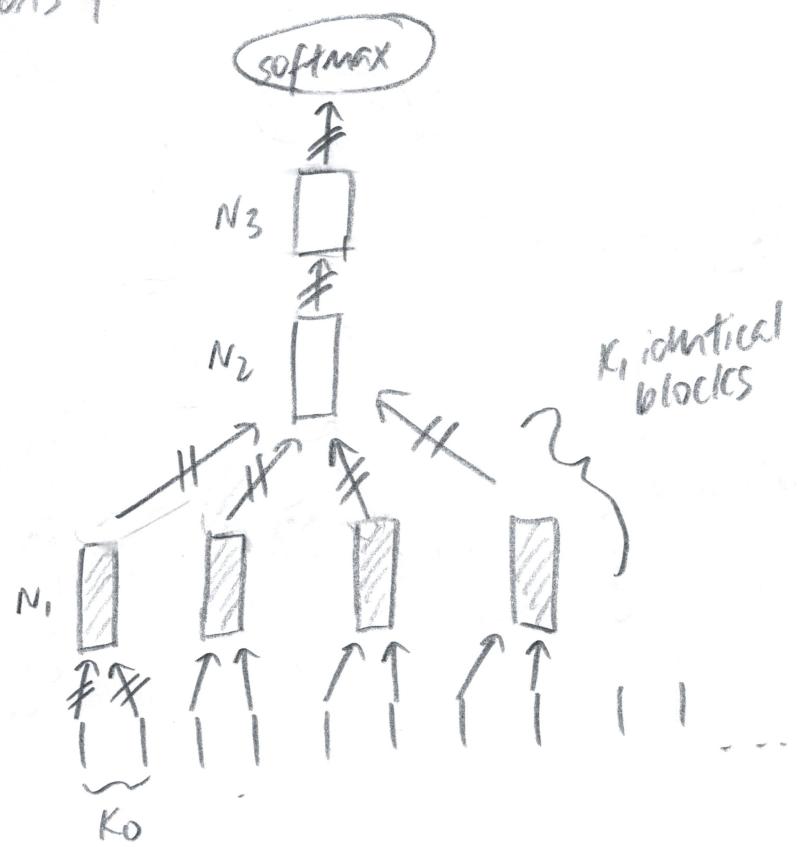
convolutional neural network (SM6)

(aka scanning without MLP)

- The key point here is the definition of convolution; or exposure to idea.
- This is covered in more detail and defined in CNN pt I and II (late lectures).
- Also note it is the re-ordered, layer first iteration computation that is basis for CNNs

- (6) that is, for precursor to CNNs,
we do ~~to~~ iterate through layers first
- the 3D tensor / 3D scan
is a bit obtuse to me (slide 197).
- ### 3) counting parameters
- scanning without distribution (slides in S.2020 suppl. slides part III)
-
- Each arrow represents edges / weights → see suppl. slides S.2020 pt(I)
- (7) Best to see this in context of how multiple copies of identical subnets (i.e. with shared parameters) feed into a final softmax layer for non-distributed scanning.
- counting parameters.
- We have a window of 8 time steps (or t time steps in general).
 - Each window 8 input vectors, each of dimensionality D
 - There are N_1 neurons in the 1st block.
 - There are N_2 neurons in the 2nd block.
 - There are N_3 neurons in the 3rd block
 - Hence total no. of params N_W :-
- $$N_W = 8DN_1 + \underbrace{N_1N_2}_{\textcircled{I}} + \underbrace{N_2N_3}_{\textcircled{II}} + \underbrace{N_3}_{\textcircled{III}} + \underbrace{\text{softmax}}_{\textcircled{IV}}$$
-

- In general, for k time instances/window width:-
 $N_w = kDN_1 + N_1N_2 + N_2N_3 + N_3$ (in a subnet)
- This ignore bias terms in the computation; this for without distribution
- We only count parameters for one column (i.e. subnet I think) since other columns/subnets are identical.
- Scanning (without distribution) uses these shared/identical parameters
 i.e. the same subnet scanning
- Distributed scanning (slides in S.2020 suppl. slides part II)
- Currently, on review, it is unclear how the partition maps to the parameter quantity computation.
- List the parameter quantity calculations for distributed and non-dist case (there are ambiguities.)
- This new diagram represents a scan of 8-time step wide patterns with a stride of 2 time steps distributed over layers.
- All double marked arrows \leftrightarrow denotes the unique weights for distributed scanning
- This is what we count.
- Lower level shaded blocks are identical i.e. clones of one another
- Counting parameters:
 - A window of 2 time steps (or K_0 timesteps in general).
 - Each



18-CNNs-review

- Supplements
key points in notes
and slides.

(*) Perceptron as a correlation filter

- Note:-

$$y = \begin{cases} 1 & \text{if } \mathbf{x}^T \mathbf{w} \geq T \\ 0 & \text{otherwise} \end{cases}$$

- we can view this through the lens of correlation

- Perceptron fires if correlation between weights and inputs exceed a threshold

1.

(*) Perceptron fires if input pattern looks like pattern of weights.

(*) Following this; can view input layer as 'feature detectors'

that detect if certain patterns have occurred in input.

(*) Network can be viewed as a function over feature detectors

(*) SHIFT-INVARIANT

(*) Scanning with an MLP.

- Sliding windows \rightarrow each produces a classification

- Combination into final output :-

- OR function: - if any window contains 'welcome' (speech recognition)

- equivalent to MAX function

- softmax/logistic: adjacent windows can combine evidence.

(*) Clarity on input of MLP: - later.

(*) Basic scanning for 1D/2D input pseudocode.

- KxK - MLP patch-size H-height of image
W - width of image (2D-image)

FOR i=1:W-K+1

FOR j=1:H-K+1

Img segment = img(i:i+W-1, j:j+H-1)

y(i,j) = MLP(ImgSegment)

y = softmax(y(1,1), ..., y(W-K+1, H-K+1))

'1D-data'

(R1) (*) See end
of these
notes.

- K-width of patch

FOR t=1:T-K+1

xsegment = x[:, t:t+K-1]

y(t) = MLP(xsegment)

y = softmax(y(1), ..., y(T-K+1))

(R2)

(*) Regular networks vs scanning networks

- Regular NN:

- each neuron in a layer is connected by a unique weight to every unit in previous layer (may be an input)

- E.g. first layer weights:-

$$\underline{W}^{(1)} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1m} \\ w_{21} & w_{22} & \dots & w_{2m} \\ \vdots & & & \\ w_{N1} & w_{N2} & \dots & w_{Nm} \end{bmatrix} \quad W \in \mathbb{R}^{N \times M}$$

mique

(*) weights matrix: $(N \times M)$ fully populated with NM parameters

- Scanning NN:

- each neuron is connected to a subset of neurons in the previous layer.
- sparse weights matrix
- weights matrix structured with identical blocks

$$\underline{W}^{(1)} = \begin{bmatrix} w_{11} & w_{12} & & & \\ w_{21} & w_{22} & & & \\ w_{31} & w_{32} & & & \\ & & w_{11} & w_{12} & \\ & & w_{21} & w_{22} & \\ & & w_{31} & w_{32} & \\ & & & & w_{11} & w_{12} \\ & & & & & w_{31} & w_{32} \end{bmatrix}$$

- network is a shared parameter network

(*) stride \rightarrow need to clarify

(*) param estimation in shared param nets:-

$$\frac{\partial \text{Div}}{\partial w^S} = \frac{\partial \text{Div}}{\partial w_{ij}^R} \frac{\partial w_{ij}^R}{\partial w^S} + \frac{\partial \text{Div}}{\partial w_{mn}^L} \frac{\partial w_{mn}^L}{\partial w^S} = \frac{\partial \text{Div}}{\partial w_{ij}^R} + \frac{\partial \text{Div}}{\partial w_{mn}^L}$$

$= 1 \qquad \qquad = 1$

- rest are zero entries

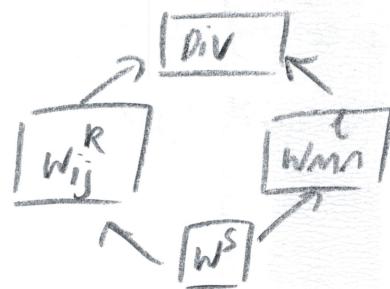
-? unclear

- especially
as there's a
time-dimension

- do we use
a weights matrix
that is a large placeholder?
i.e. a $K^2(NM)$ as logic
where K is patch size?

as $w^S = w_{ij}^{(R)} = w_{mn}^{(L)}$

Influence diagram:



(*) computing divergence of shared parameters.

- define $S = \{e_1, \dots, e_n\}$
- any set of edges that have a common value, and w^S is the common weight of the set. (see slides in colow)

$$(1) \frac{\partial \text{Div}}{\partial w^S} = \sum_{e \in S} \frac{\partial \text{Div}}{\partial w^e}$$

(*) individual terms in sum \rightarrow use backprop

(*) training networks with shared param. PSEUDO CODE

- gradient descent algorithm
- initialise all weights w_1, w_2, \dots, w_R

DO:-

- for every set S

$$\nabla_S \text{Err} = \frac{\partial \text{Err}}{\partial w^S}$$

$$w^S = w^S - \eta \nabla_S \text{Err}$$

- for every $(k, i, j) \in S$ update:-

$$w_{i,j}^{(k)} = w^S$$

- until Err has converged

For every training instance X

- for every set S :-

$$\nabla_S \text{Div} \dagger = \frac{\partial \text{Div}}{\partial w_{i,j}^{(k)}} \quad (*)$$

$$\nabla_S \text{Err} + \nabla_S \text{Div} \dagger$$

(*) computed via backprop

(*) scanning order (clarity)

1) Pass each window through entire network, this happens N times

for each window w_i

- these N outputs are then passed through softmax

2) Pass each window through neuron, neuron by neuron, layer by layer.

- 1st neuron - output for W windows 1st layer

- 2nd neuron " "

⋮

- 1st neuron " " 2nd layer

⋮

(*) (*) order of computation merely amounts to exchanging loop nesting order
 ↳ does not change final output

(*) scanning order pseudocode (1-D)

order 1 - Window by window through entire net. i.e.

FOR $t = 1:T-K+1$ → over window

→ over layers

FOR $l = 1:L$

FOR $j = 1:D_l$

if ($l=1$) # first layer operates on input

$y[0,:,t] = x[:,t:t+K-1]$

end

$z(l,j,t) = 0$

FOR $i = 1:D_{l-1}$

$z(l,j,t) += w(l,i,j) y(l-1,i,t)$

$y(l,i,t) = \text{activation}(z(l,i,t))$

$y = \text{softmax}(y(l_1,:), \dots, y(l_i, T-K+1))$

(*) exchanged order 2:-

FOR $l = 1:L$

FOR $t = 1:T-K+1$

FOR $j = 1:D_l$

(*) Scanning MLP-vector not. (1-D)

FOR $i = 1:L$

FOR $t = 1:T-K+1$

if ($t == 1$)

$$y(0,t) = x[1:t+K-1]$$

end

$$\underline{z}(l,t) = w(l)y(t-1,t)$$

$$y(l,t) = \text{activation}(\underline{z}(l,t))$$

$$y = \text{softmax}(y(l,1), \dots, y(l,T-K+1))$$

(*) Can view activity at one neuron as:-

$$\left(\sum_{i,j} w_{ij} p_{ij} + b \right) \quad \begin{matrix} \text{for one window} \\ \text{of input data } p \end{matrix}$$

(*) (67): view slides interactively \rightarrow helps with the visual intuition of slides (excellent on GR part)

(*) Scanning - a closer look

- recursing the logic

(*) Second-level neurons are scanning the rectangular outputs of the

1st level neurons.

(*) They are jointly scanning multiple 'pictures'

(*) Each location in output of 2nd level neuron considers corre. locations from outputs of all first-level neurons

(*) Output of penultimate layer is a grid/picture

(*) This entire grid is passed through a final layer/neuron that performs:-

- i) logical OR
- ii) max output
- iii) softmax
- iv) MLP

(*) 2-D scanning pseudocode

(*) Non vector form

FOR $x = 1 : W - K + 1$

FOR $y = 1 : H - K + 1$

FOR $l = 1 : L$

FOR $j = 1 : D_l$

IF ($l == 1$) # first layer ops on input

$y(0, :, x, y) = \text{img}(1 : C, x : x + K - 1, y : y + K - 1)$

end

$z(l, j, x, y) = 0$

FOR $i = 1 : D_{l-1}$

$z(l, j, x, y) += w(l, i, j) y(l-1, i, x, y)$

$y(l, j, x, y) = \text{activation}(z(l, j, x, y))$

$y = \text{softmax}(y(l, :, 1, 1) \dots y(l, :, W - K + 1, H - K + 1))$

(*) Adjusted order:

FOR $l = 1 : L$

FOR $j = 1 : D_l$

FOR $x = 1 : W - K + 1$

FOR $y = 1 : H - K + 1$

(*) Vector version \rightarrow see slides.

(*) Distributing the scan/behaviour of layers.

① notes were a little opaque on this

②/5 ③: finding the distinction between earlier presentation and explanation of how scan is distributed difficult to follow

(*) In this context:

- stride - no. of steps we move 'sliding window' across by
- also have to specify size of window

(R1): Here, we have n inputs $x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(n)}$ (1-D case)

• each i^{th} input is a $0 \times T$ matrix.

• D -dimensions and T -temporal length.

• Notes from this point onwards clarify things which come up.

$$x^{(i)} = \begin{bmatrix} x_{01} & x_{02} & \dots & x_{0T} \\ \vdots & \vdots & \ddots & \vdots \\ x_{01} & x_{02} & \dots & x_{0T} \end{bmatrix}$$

• for 10 sequences; if we want a width K patch to scan over $x^{(i)}$, then we do

$x^{(i)}[:, t:t+K-1]$. → slices are using Python slice notation.

• for clarity, for a window of width K :-

$t, t+1, t+2, \dots, t+K-1, t+K$

a window of width K (i.e. accounting for slicing notation).

• for clarity on the loop

FOR $t=1:T-K+1$

$x\text{-segment} = x[:, t:t+K-1]$

$y(t) = \text{MLP}(x\text{-segment})$

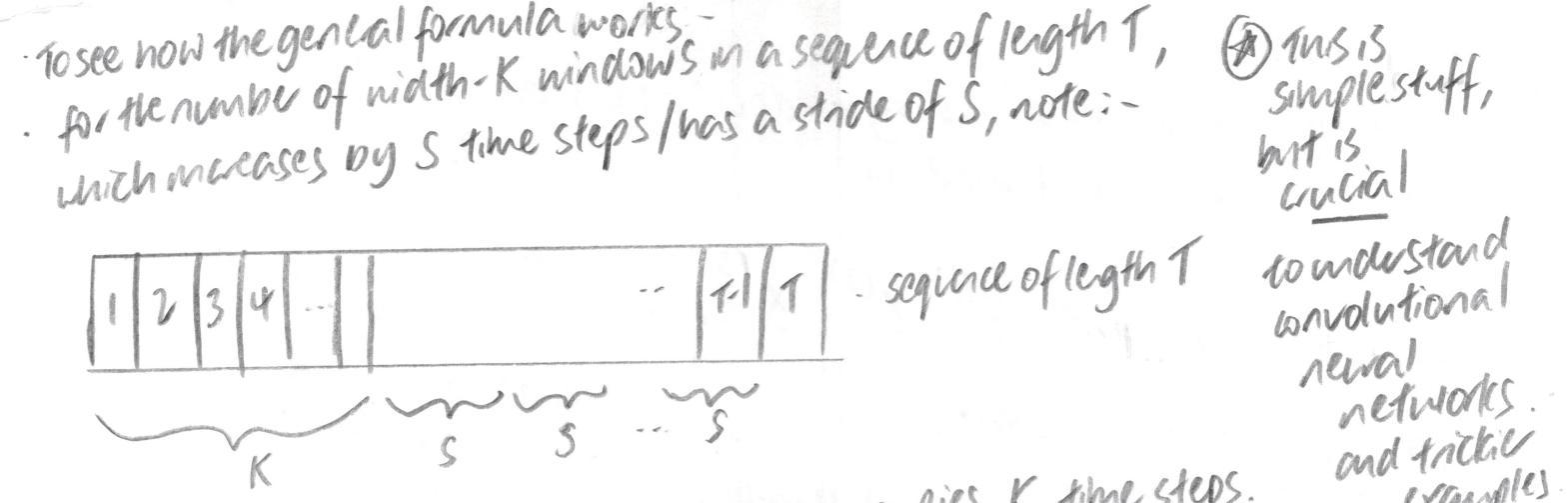
$y = \text{softmax}(y(1), \dots, y(T-K+1))$

• At $t=1$
 $x\text{-segment} = x[:, 1:K]$

• At $t=T-K+1$
 $x\text{-segment} = x[:, T-K+1:T]$

• we are proceeding along a length T sequence K steps at a time.

• It is important to note that we proceed with a width K sliding window and the no. of steps this sliding window advances by 1 time step at each point. That is the above example has stride $S=1$.



- In its initial position, the width-K window occupies K time steps.
- It will advance with a stride length of S time steps, until termination.
- We want to know the number of width-K windows, which we call N.
- Notice that if we subtract K time-steps from the total length of the sequence T, we have $(T-K)$.
- Then the no. of strides of length S taken by the sliding window is just :-

$$\frac{T-K}{S}$$

- Accounting for the fact that $(T-K)$ may not be a multiple of the stride length S, we carry out remainder division, and forget about remaining steps at the end.
- The number of remaining time steps at the end will be less than S (otherwise we could allow sliding window to progress another S timesteps).
- Hence the no. of strides of length S taken by the sliding window is :-

where $\lfloor \cdot \rfloor$ is the floor function

$$\left\lfloor \frac{T-K}{S} \right\rfloor$$

- As the number of width-K windows is just one more than the number of strides, i.e. accounting for the initial position of the sliding window, the general formula is

$$N = \left\lfloor \frac{T-K}{S} \right\rfloor + 1$$

- In the case where stride length $S=1$ (i.e. in the previous example on the slides)
 - e.g. $y = \text{softmax}(y(1), \dots, y(T-K+1))$ earlier
 - Note that here $S=1$, and the formula for no. of sliding windows is

$$N = \left\lfloor \frac{T-K}{1} \right\rfloor + 1 = T-K+1 \quad \text{as required.}$$
-
- (R2) - On the dimensionality of:-
- $y = \text{softmax}(y(1), y(2), \dots, y(T-K+1))$.
- $Y \in \mathbb{R}^{M \times (T-K+1)}$ in the time-delay NN/ID example,
where M is no. of classification categories (if we're doing classification).
 - in the example in slides of detecting whether there is "welcome" in
an audio signal.
 - note that we could model with $M=2$ i.e. yes or no
 - The softmax output $y(i) \in \mathbb{R}^2$ in this case (consisting of probabilities).
 - The reason for Y having $(T-K+1)$ columns is because we want probabilities
over 2 classes for each of $(T-K+1)$ individual windows of width K .
(or. Softmax
stage)

