

18-CNNs

(*) Next 4 lectures → scanning for patterns
(2 weeks) (CNNs)

(*) Scott Fahlman - guest lectures

(*) New problem - spectrogram

(*) From this data, can we check if the signal contains 'welcome'

BR: Yes, as an I/O problem

(*) Issue: - Two signals have welcome stated at diff. locations

(*) Q: will MLP trained on signal with 'welcome' give the same classification if the signal with 'welcome' is shifted in time?

- NO, covariates have different values; no guarantee.

- consider through vector representations

BR: we want network to fire, regardless of location → shift/translation invariance

(*) Visual analogies

- A flower 'shifted' within a picture
- patterns lie on different subspaces

(*) we require shift invariance for both speech and visual recog.

(*) Solution: - use scanning i.e. a sliding window over inputs

and classify each window.

How to combine classification over each window frame to give a classification

- use max or or or logistic

(*) BR: Switched issue of whether recording contains word 'welcome'
to one of scanning and then combining outputs

Scanning with an MLP

(*) A1 - see pseudocode ✓

- can view as one large neural network
- NN composed of many subnets (same one)
- subnets are identical
- This is a shared parameter network

(*) Scanning to find a flower

- we have a classification for each window of the picture
- pass all of these interim class through a max/softmax.

(*) Scanning image pseudocode.

(*) A2 - Go through. Subnet 1

- for MLP, not just scan $L \rightarrow K$ but $T \rightarrow B$.

(*) Network training

- constant of shared parameters → each ^{lower-level} subnet that classifies windows frame are identical

(*) Any update of param of one copy → update all copies.

(*) Learning in shared param networks

- shared params → arbitrarily deciding that some params. have same value

(*) e.g. $w_{ij}^R = w_{mn}^l = w^S$

(*) for training instance X , small perturb. of w^S perturbs both
 w_{ij}^R and w_{mn}^l identically

- each perturbation will individually influence $\text{div}(d, y)$

• Attempting to estimate common value w_s

(*) computing divergence of shared params.

(a3)-review influence diagram
of Div $\rightarrow w_s$

(*) computing divergence of shared params

(a4)-review - formula

(*) Training a network with shared params.

(a5)- pseudocode + detail - review

- slight modification on backprop (76 slides in 18 mins)

Story so far

(*) position-invariant pattern class \rightarrow scanning

(*) Scanning \Leftrightarrow composing a large network with repeating subnets.

(*) Param estimation \rightarrow backprop rules must be modified to combine gradients from param that share same value. (applies to all networks with shared param)

(*) Scanning - closer look.

- visual illustration \rightarrow very clear

- order of computation / operations. (scanning)

- for W windows of spectrogram.

{ - W_1 - window 1 - pass through 4 unit, 2 unit, 1 unit layers

{ - W_2 - 11 - 2

{ - W_W - window W

\hookrightarrow pass all through softmax

④ Window by window

- (*) different order
- pass each window through individual unit (by unit), layer by layer.
 - does not change output
- Ⓐ - seti Review. ✓

- (*) Same computation
- columnwise vs row wise - both scanning MLPs
 - original pseudocode
 - scan L→R, pull out segment, MLP it; complete softmax
 - different order pseudocode.
 - run through neurons
- Ⓐ psendocode
- (*) can switch for loop nesting order; does not change comp.
- (*) scanning with MLP- vector notation → Ⓑ Ⓕ /

- (*) Scanning flower image
- 1) window by window (entire MLP)
 - 2) neuron by neuron (scan entire image with one neuron)
 - 3) output for that neuron for every window of image. (scan image)
 - do some for every neuron
 - we then have a map of the image at every neuron
 - (i.e. a collection of outputs over windows at every neuron)

- (*) To find whether there is a flower / classify a patch
- Ⓑ Pick those locations from it maps; pass it through MLP.
- Ⓒ: note two equivalent ways of conducting, operations ^{scan}
- (*) process all neurons in one location
- (*) process all locations with each neuron
- (*) logic can be reused at next layer.

(*) Scanning - 2nd layer neurons / 3rd output?

- At next level
- scan outputs of 1st layer
- (*) To find if there is a flower in a particular location, pick the location from the maps ...
- (*) review logic - crucial ⑦

(*) Redefining final layer?

- still scanning MLP, except change in comput. order; but still scanning MLP with shared params.

- ⑧ changes in pseudocode, comparison of window by window, neuron by neuron

OR: forward idea of MLPs / NNs capturing progressively more abstract/complex patterns hierarchically

(*) Distributing the scan/behavior of the layers.

- 1st input \rightarrow pixel by pixel (?)

(*) First layer neurons resp. for evaluating entire block of pixels. (2) ⑨ ⑩ ⑪ - logic

- subsequent layers look at single pixel on input maps

(*) Distribute scanning over 2 layers

(to achieve some block analysis at 2nd layer)

- first layer \rightarrow smaller block of pixels

- next layer \rightarrow blocks of output from first layer

- effectively evaluates large block of original image

(*) still scanning with a shared param network

- with minor modification. \rightarrow ⑪ - clarity or parameter sharing

- (A11) - smearing information over layers
- recursively decomposing patterns over layers

(*) Hierarchical build up of features.

- scanning and combining outputs as giant shared MLP
- individual nets have shared param sets \rightarrow does not change.

- (A12) pseudocode.

- can re-order computation again. \rightarrow but size of maps change. (A13)
pseudocode

- Affine combination \rightarrow review details
- inner product between weights and inputs with transl.

(*) (A14) - This is a convolution \rightarrow convolutional neural network

- NN with scanning
- typically, scan input at finer and finer resolutions as we go through the network.

(*) CNN - vector notation (A15)

(*) Why distribute?

(D) distribution forces localised patterns in lower layers

(D) (why not one single scanning MLP that scans entire image?)

(*) scanning without distri

- need to print new slides S.2020 significantly different

(*) why distribute?

- distib. - localized patterns in lower layers
- no. parameters - large reduction (OOM) reduction in param.
 - significant shared computation gains.
- Regardless of how scanning is distri. computationally ; still a scanning MLP.

(*) story so far

(*) operations in scanning input with full network can be reordered

- scanning the input with individual neurons in the first layer to produce scanned 'maps' of input
- jointly scanning the 'map' of outputs by all neurons in previous layers by neurons in subseq. layers

(*) Scanning block can be distributed over multiple layers of network

- significant reduction in total no. param.

(*) Hierarchical composition - different resp.

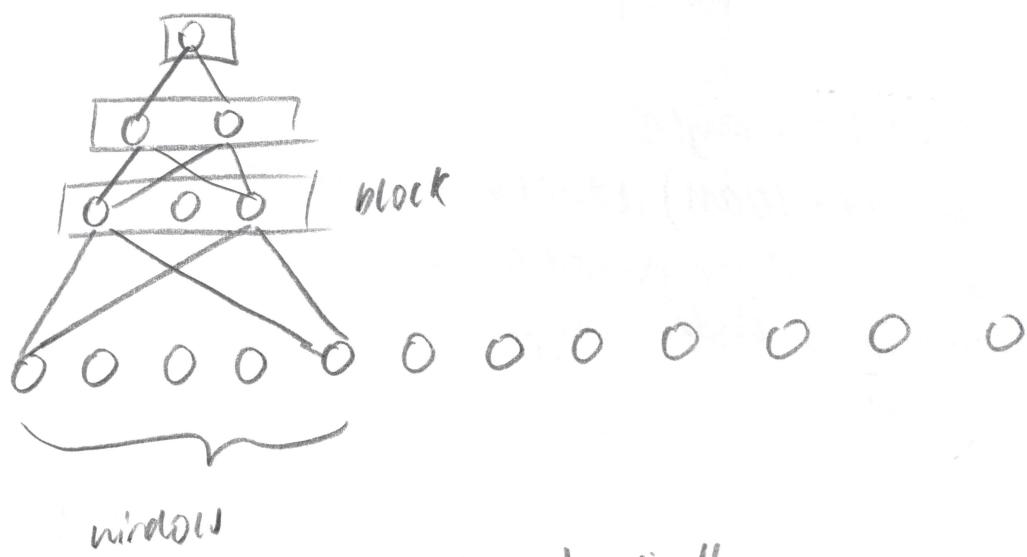
- 1) first layer \rightarrow subregions of math image (eg. petals)
- 2) 2nd layer \rightarrow regions of output of first layer.
 - int petals \rightarrow flower
 - large region of origin. input image

(*) Each of the scanning neurons as a filter

- A correlation filter
- Each filter scans for a pattern in the map it operates on

(*) Pattern in the input image each filter sees \rightarrow receptive field

- first layer filter: actual receptive field is simply arrangement of weights
- higher level filters: - receptive field not obvious, must be calculated.
 - may not have natural correspondences due to complex combos.



- Rearrange neurons in a block (layer) vertically

- $N_1 \rightarrow N_2$ - 6 connections (N_1, N_2)

- $N_2 \rightarrow N_3$ - 2 connections (N_2, N_3)

(*) → check/review $\xrightarrow{\text{calculation}}$ formulae or no. param for generic scanning vs distributed scanning. (parameter savings) calc.



N_1 neurons

0 0 0 0 0 0 0 0

- yellow columns

- snapshot

- columns

of N_1 neurons
scanning 2/k₀
patterns

at a time.

(*) Distributed scanning \Rightarrow few parameters to estimate

(*) large additional benefit
from the fact that scans at neighboring positions share the computation
of low-level blocks

(*) note unique parameters

(*) re-use of computation \rightarrow more efficient

(*) go through 2D case in your review

(*) - 1034 parameters vs 160 parameters
(non-distri. scanning) (distrib. scanning)

(*) Modifications

- final layer may feed into MLP rather than single neuron.
- shared param net.

(*) mod 1 - conv 'Stride'

- stride is key concept (for CNNs)
- Pseudocode A18

(*) Accounting for jitter

- minor noise

(*) Similar principle

- Is there a flower in image?
 - scan entire input, take a max of all outputs \rightarrow tells us whether there is a flower or not

(*) Small jitter acceptable

- replace each value by maximum of values within a small region around it
- max filtering/pooling
- Max is just an activation

B) equivalent to adding a layer with weights all one; activation is max. that follows previous layer
B) does not change network as shared param.

(*) Max pooling 'strides'

(*) Max operations may stride by more than one pixel

- shrinks map
- operation is called 'pooling'
- consists of pooling a no. of outputs to get a single output.
- $\text{stride} > 1 \rightarrow \text{downsampling}$

(*) next layer works on max-pooled maps

①: overall structure

- many layers of convolution (scanning) followed by max pooling (and reduction) before final MLP.
- (*) individual perceptions at any scanning/convolutional layers are called filters.
- (*) filter image (input) to produce output image (map)
- (*) Indiv max operations → max pooling filters
- (*) entire structure → CNN
- scanning for patterns, distributing manner in which they are represented across network.

(*) shift invariance

- LeCun
- MNIST digit scanning (not on individual frames)

(*) 1D convolution

- speech
- still convolution (but across time)
- representation:- spectrographic time frequency components
- no jitter (and max pooling) tolerance on speech

(*) story so far - summary