

convolution
(filter x features)

→ many useful image features are local

↳ apply a "mini-layer" (aka filter) to process patches of image.

↳ "downsample" activation grid after filtering (aka pooling).

→ typically: Max Pooling 2D.

→ historical example: LeNet, handwritten digit recognition.

→ a way of avoiding needing millions of parameters w/ images
each layer is "local", produces an "image" w/ (roughly) the same width and height. # of channels = # of filters.

N-dim (ND) arrays \leftrightarrow tensors

this is the feature @ i, m channel

Conv: (F_H, F_W, C_H, C_W)

$$z^{(2)}[i, j, k] = \sum_{l=0}^{H_F-1} \sum_{m=0}^{H_W-1} \sum_{n=0}^{C_H-1} W^{(2)}[l, m, k, n] a^{(1)}[i+l+\frac{H_F-1}{2}, j+m-\frac{H_W-1}{2}, n]$$

$$z^{(2)}[i, j] = \sum_{l=0}^{H_F-1} \sum_{m=0}^{H_W-1} W^{(2)}[l, m] a^{(1)}[i+l-\frac{H_F-1}{2}, j+m-\frac{H_W-1}{2}]$$

$a^{(2)}[i, j, k] = \sigma(z^{(2)}[i, j, k])$ apply non-linearity.

$$H_{out} = H_{in} - \frac{H_F-1}{2} \times 2 = H_{in} - H_F + 1$$

Input: $32 \times 32 \times 3$
Filter: $5 \times 5 \times 6$ Output: $28 \times 28 \times 6$

A way to tackle padding w/o losing dimensionality: zero-pad

Strided convolutions: skip over some positions w/ our filter. Saves computation.
→ ppl argue, just as good as conv + pooling.

Examples:

→ AlexNet, obtained SOTA results on ILSVRC (work of art), 8 layers
trained on two GPUs!

Conv 1: $11 \times 11 \times 96$, stride 4, maps $224 \times 224 \times 3 \rightarrow 55 \times 55 \times 96$ (w/ zero pad)

Pool 1: $3 \times 3 \times 96$, stride 2, maps $55 \times 55 \times 96 \rightarrow 27 \times 27 \times 96$

Norm 1: local normalization layer (not widely used anymore)

Conv 2: $5 \times 5 \times 256$, stride 1, maps $27 \times 27 \times 96 \rightarrow 27 \times 27 \times 256$ (w/ zero pad)

Pool 2: $3 \times 3 \times 256$, stride 2, maps $27 \times 27 \times 256 \rightarrow 13 \times 13 \times 256$

Norm 2: local norm layer

Conv 3: $3 \times 3 \times 384$, stride 1, maps $13 \times 13 \times 256 \rightarrow 13 \times 13 \times 384$ (w/ zero pad)

Conv 4: $3 \times 3 \times 384$, stride 1, maps $13 \times 13 \times 384 \rightarrow 13 \times 13 \times 384$ (w/ zero pad)

Conv 5: $3 \times 3 \times 256$, stride 1, maps $13 \times 13 \times 384 \rightarrow 13 \times 13 \times 256$ (w/ zero pad)

Pool 3: $3 \times 3 \times 256$, stride 2, maps $13 \times 13 \times 256 \rightarrow 6 \times 6 \times 256$

FC 6: $6 \times 6 \times 256 \rightarrow 9216 \rightarrow 4096$ [matrix is 4096×9216]

FC 7: $4096 \rightarrow 4096$

FC 8: $4096 \rightarrow 1000$

Softmax

→ VGG (work of engineering), still often used today, homogeneous stacks of convs.

↳ more layers, more processing, bulk of parameters in FC 1...

→ ResNet (152 layers!), longer stacks of convolutions

↳ no large FC layer at end, just average pool last conv, results in one linear layer, reduces

↳ implements residual (skip) connections: $H(x) = F(x) + x$ (each input to output) parameters \leftarrow 2 layers
(good gradient flow)

more features,
more things
could be
"going on"

tradeoff,
features, res1
ReLU every
conv/FC layer
but last

helps stabilize
dx
neither 0
nor too

19 layers
↓

CS 182 Lecture 7: Initialization, BatchNorm

6/15

$$\frac{dL}{dw^{(l)}} = \frac{dz^{(l)}}{dw^{(l)}} \frac{dL}{dz^{(l)}} = Sx^T \leftarrow \text{if } x \text{ is "imbalanced", so are gradients.}$$

Therefore, we really want all entries in $x \rightarrow$ same scale.

standardization: $\mu=0, \sigma=1$.

$$\begin{aligned} \bar{x}_i &= x_i - E[x] \\ \bar{x}_i &= \frac{x_i - E[x]}{\sqrt{E[(x_i - E[x])^2]}} \leftarrow \bar{\sigma}_i \end{aligned}$$

$E[x] \approx \frac{1}{N} \sum_{i=1}^N x_i$

Basic idea of Batch Norm: standardize activations at each layer, controlling gradients. computationally expensive idea, so only perform on batch.

$$\mu^{(l)} \approx \frac{1}{B} \sum_{i=1}^B a_i^{(l)} \quad \sigma^{(l)} \approx \sqrt{\frac{1}{B} \sum_{i=1}^B (a_i^{(l)} - \mu^{(l)})^2} \quad \bar{a}_i^{(l)} = \frac{a_i^{(l)} - \mu^{(l)}}{\sigma^{(l)}} \gamma + \beta$$

\leftarrow learnable scale and bias same dim as $\bar{a}_i^{(l)}$

- \rightarrow Can be trained via backprop, since these are differentiable.
- \rightarrow Can be placed either before/after nonlinearity.
- \rightarrow We can often use a larger learning rate
- \rightarrow models can train faster
- \rightarrow generally requires less regularization

Basic initialization methods: ensure activations are on reasonable scale that stay constant. More advanced init. methods involve eigenvalues/Jacobians. "Try" to have well-behaved gradients from the get-go.

Set $W_{ij} \sim N(0, \sigma^2_w)$, $b_i \approx 0$, $z_i = \sum_j W_{ij} a_j + b_i$ assume $a_j \sim N(0, \sigma^2_a)$

$$\text{so, set std of } W_{ij} = \frac{1}{\sqrt{D_a}} \quad E[z_i^2] = \sum_j E[W_{ij}^2] E[a_j^2] = D_a \sigma^2_w \sigma^2_a$$

\leftarrow dim of a

"Xavier initialization" if $D_a \sigma^2_w \gg 1$, magnitude grows w/ each layer

ReLU issue: zero out activations. If $D_a \sigma^2_w < 1$, magnitude shrinks w/ each layer $\rightarrow \sigma^2_w = \frac{1}{D_a}$

Ergo, increase std of $W_{ij} \rightarrow \frac{1}{\sqrt{2} D_a}$ (proposed in ResNet).

Alternative, init. $b_i = 0.1$ (or small constant) to avoid zero-out in ReLU.

Reminder: $\frac{dL}{dw^{(l)}} = J_1 J_2 J_3 \dots J_n \frac{dL}{dz^{(n)}} \quad \forall x, J_i = U_i \Lambda_i V_i^T$ (SVD). diag. matrix of λ 's of J_i .

By ensuring λ 's are scaled, we avoid 0/infinity convergence issue.

Since $J = W^T$ (derivative w) $\rightarrow W^{(l)} = U^{(l)} \Lambda^{(l)} V^{(l)T}$, force $\Lambda^{(l)} = I$, $W^{(l)} \approx U^{(l)} V^{(l)T}$

"Measure of last resort": Gradient clipping, because "monster gradients" can occur

\rightarrow per-element clipping: $\bar{g}_i \leftarrow \max(\min(g_i, c_i), -c_i)$

\rightarrow norm clipping:

$$\bar{g} \leftarrow g \frac{\min(\|g\|, c)}{\|g\|} \leftarrow \text{clips length not dimension, see what "healthy" magnitudes look like (plot), choose } c \text{ thru experimentation}$$

NNS \rightarrow high-variance (lots of parameters), but with multiple, more agreement on "right"

$$\text{Variance} = E_p(p) [\|f_{\theta}(x) - f(x)\|^2] \quad f(x) = E_p(p) [f_{\theta}(x)] \approx \frac{1}{M} \sum_{i=1}^M f_{\theta_i}(x) \quad \text{bootstrap sample}$$

Select? \rightarrow principled approach: average $(p_{\text{avg}}(x) = \frac{1}{M} \sum_{i=1}^M p_{\theta_i}(x))$, often more robust

\rightarrow simple approach: majority vote (analogous to "first past the post")

Ensembles in practice: train M models $p_{\theta_i}(x)$ on the same \mathcal{D} + training set, maj. vote. (even faster ensembles? share features for all M models. Separate ensemble of fixed classifier "heads" at end (then are task-specific))

Snapshot ensemble: save our parameters as "snapshots" during SGD, use later as a model. The bigger the ensemble, the better, the more expensive.

Dropout: randomly set some activations to 0 in forward pass. Creates a "new" network, helps w/ ensembling out of a single NN. 2^{km} diff. models (k layers, m nodes/layer)

At test time $\bar{W}^{(l)} = \frac{1}{2} W^{(l)}$, since on average, $\frac{1}{2}$ of dimensions are forced to 0.

Hyperparameters that affect generalization (validation): ensembling, dropout, arch.

General method to pick hyperparameters: coarse-to-fine, broad sweep before "zeroing in"

\hookrightarrow In practice: random hp. search? grid search.