

Today: convolutional neural nets
(continued)

Reading: Prince through Ch 9
(so far)

Ch 108 II now.

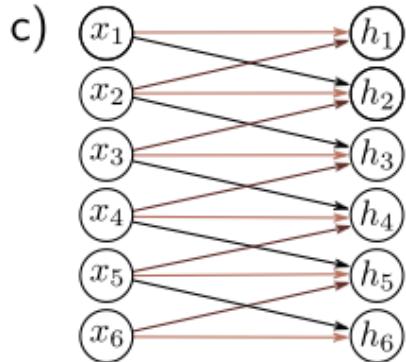
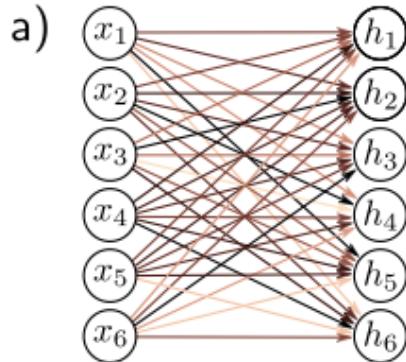
Reminder: TA OH Today

Discussion Tomorrow

Architecture Order In Class:

MLPs \rightarrow CNNs \rightarrow Graph NN \rightarrow RNN/State-space \rightarrow Transformers

We are here



b)

	x_1	x_2	x_3	x_4	x_5	x_6
h_1	Dark Brown	Light Brown	Medium Brown	Dark Brown	Light Brown	Medium Brown
h_2	Medium Brown	Black	Light Brown	Medium Brown	Dark Brown	Light Brown
h_3	Light Brown	Medium Brown	Dark Brown	Light Brown	Medium Brown	Dark Brown
h_4	Light Brown	Black	White	Dark Brown	Light Brown	Medium Brown
h_5	Black	Medium Brown	Dark Brown	Light Brown	Medium Brown	Dark Brown
h_6	Light Brown	Dark Brown	Medium Brown	Dark Brown	Light Brown	Medium Brown

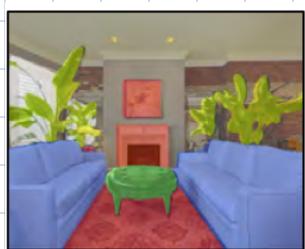
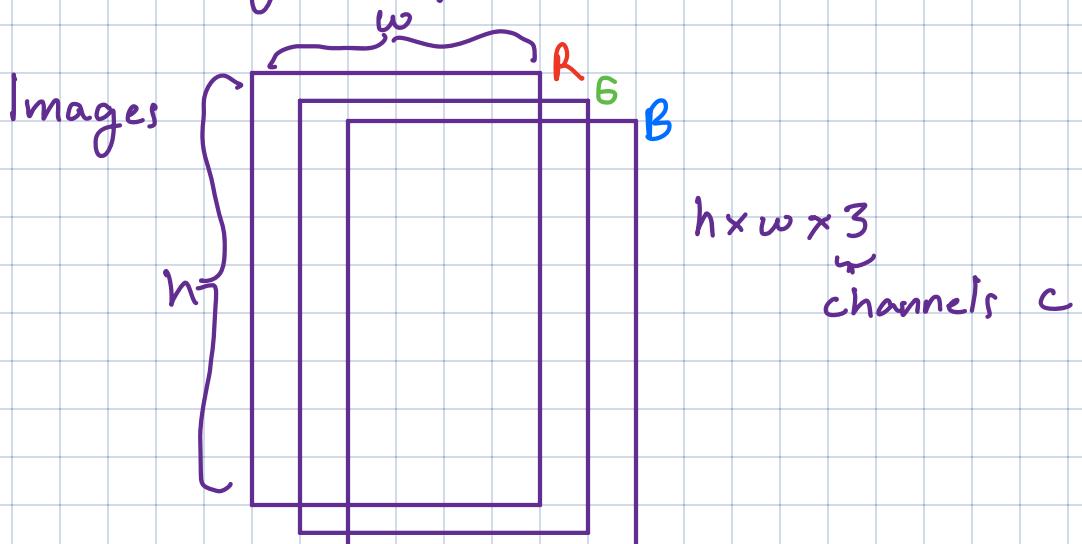
d)

	x_1	x_2	x_3	x_4	x_5	x_6
h_1	Dark Brown	Light Brown	Medium Brown	Dark Brown	Light Brown	Medium Brown
h_2	Medium Brown	Black	Light Brown	Medium Brown	Dark Brown	Light Brown
h_3	Light Brown	Medium Brown	Dark Brown	Light Brown	Medium Brown	Dark Brown
h_4	Light Brown	Black	White	Dark Brown	Light Brown	Medium Brown
h_5	Black	Medium Brown	Dark Brown	Light Brown	Medium Brown	Dark Brown
h_6	Light Brown	Dark Brown	Medium Brown	Dark Brown	Light Brown	Medium Brown

MLP

CNN

Inspired by computer vision problems: classification,



semantic segmentation,
etc.

Key Ideas (Expressivity)

- 1) Respect locality → convolutional structure with small filters.
- 2) Respect certain invariances / equivariances / symmetries → weight sharing
→ data augmentations ↑
training
- 3) Hierarchical Structure (Parts make wholes) → Depth
→ Multiresolution & Filterbanks

Key Ideas (Getting it to work)

- 1) Normalization Layers
- 2) Dropout
- 3) Residual / Skip Connections

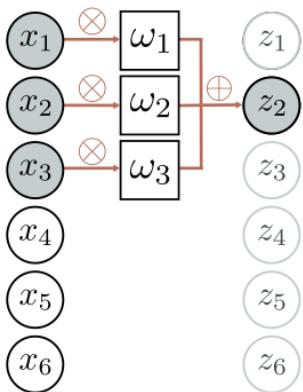
Recall 1D convolution (LTI systems... e.g. from EECS 170, 123)

$$y(t) = \int_{-\infty}^{+\infty} x(\tau) h(t-\tau) d\tau = \int_{-\infty}^{+\infty} h(\tau) x(t-\tau) d\tau$$

$$y[t] = \sum_{\tau=-\infty}^{+\infty} x[\tau] h[t-\tau] = \sum_{\tau=-\infty}^{+\infty} h[\tau] x[t-\tau] \xrightarrow{\text{FIR}} \sum_{\tau=-k}^{+k} h[\tau] x[t-\tau]$$

In Deep Learning, we don't "flip"

a)



b)

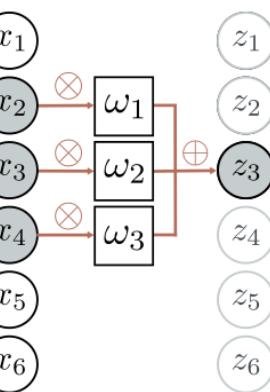


Figure 10.7 in Prince.

Notice: No Flip

Can also have a bias

a)

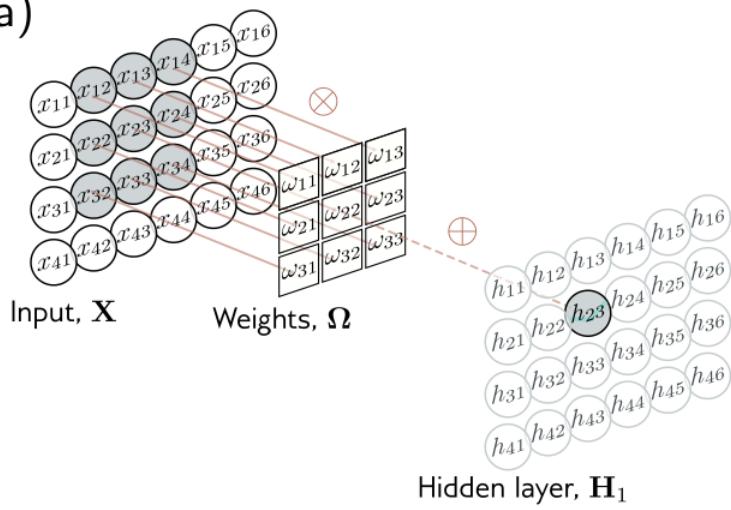


Figure 10.9 in Prince

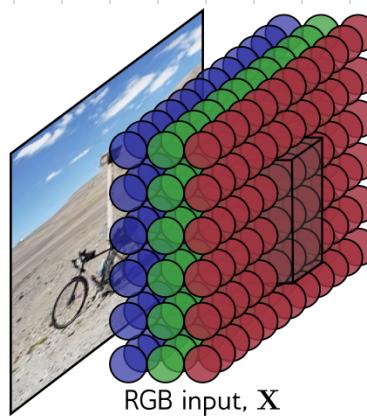
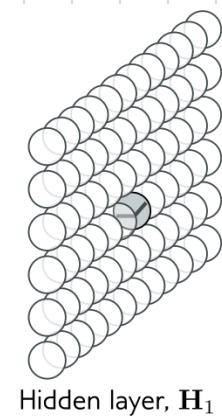


Figure 10.10
in Prince

Note: weights NOT shared across channels.



New details

Notation note: Convolutional filters called "kernels"

Examples:

1x1 conv | channel input | channel output

For every pixel, scale by a **constant** and add a **bias**
How many weights: 1
learnable weights

" " biases: 1

3x3 conv with 1 channel input 1 channel output

For every pixel, scale the value at this pixel and immediate neighbors by **individual weights**, add them up, and add **bias**.

How many weights: 9

" " biases: 1

w_{nw}	w_N	w_{NE}
w_w	w_{ne}	w_E
w_{sw}	w_S	w_{SE}

3x3 conv with C_{in} -channel input and 1 output Filter

Repeat for each channel (with individual weights per channel) and fold together by adding up.

How many weights: $9 C_{in}$

" " biases: 1

$k \times k$ conv with C_{in} -channel input and 1 output

How many weights: $k^2 C_{in}$

" " biases: 1

1×1 conv with C_{in} -channel input and C_{out} output)

Input is a vector \vec{x} } C_{in}
(from one pixel)

Output is a vector \vec{y} } C_{out}
(for this pixel)

$$\vec{y} = {}_{C_{in}}\{\vec{w}\} \vec{x} + \vec{b} \} C_{out}$$

How many weights: $C_{in} C_{out}$

" " biases: C_{out}

$k \times k$ conv with C_{in} -channel input and C_{out} output)

How many weights: $k^2 C_{in} C_{out}$

" " biases: C_{out}

Backprop?? We share parameters across different pixels

How can we compute $\frac{\partial L}{\partial w}$ when the same parameter w appears in different places in computation graph?

Easy Answer: Don't worry, PyTorch will handle it. 😊

Underlying Answer: 1) Pretend the weight w was actually $\tilde{w}_1, \tilde{w}_2, \dots$ in all the different places and just happened to be that $\tilde{w}_1 = \tilde{w}_2 = \tilde{w}_3 = \dots = w$

2) Now to get $\frac{\partial L}{\partial w} = \sum_i \frac{\partial L}{\partial \tilde{w}_i}$
i.e. Add them up

Question: How does the scale of $\frac{\partial L}{\partial w}$ change with the number of terms in Σ ?

Pooling

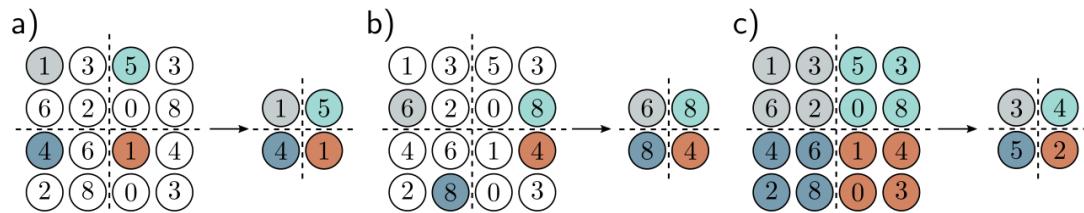


Fig 10.11 in Prince

stride/subsample

Max-Pool

Mean-Pool

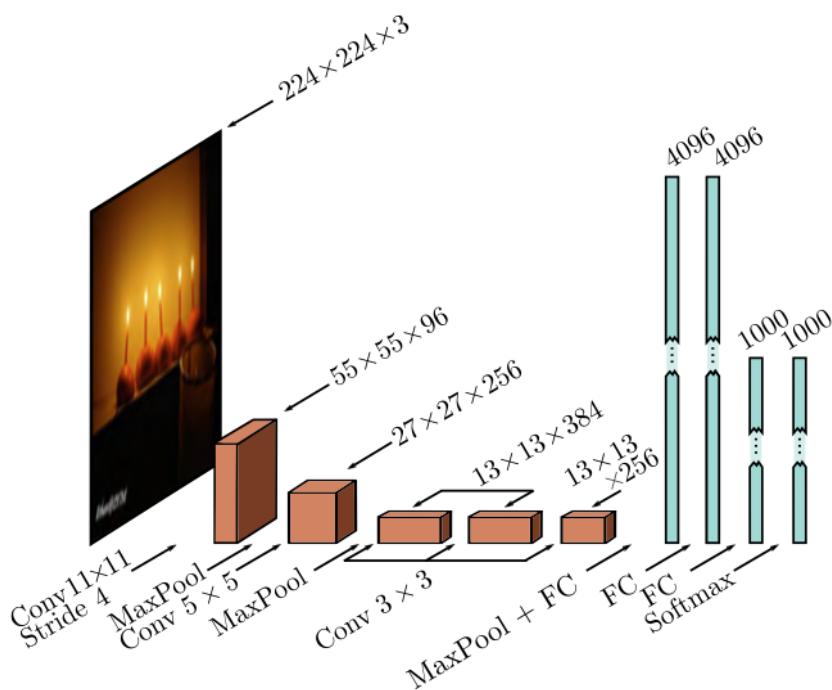
$$\text{Output} = \max_{i \in \text{Pooling domain}} h_i$$

$$\text{Output} = \frac{1}{|\text{Pooling domain}|} \sum_{i \in \text{Pooling domain}} h_i$$

Backprop? → Follow argmax

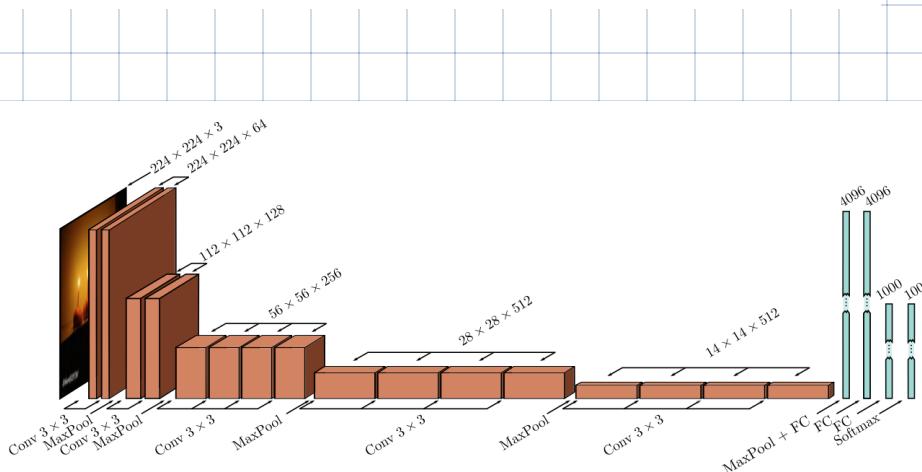
*MaxPool
(Router Gradients)*

Putting things together



Alex Net (Krizhevsky et al.) 2012

Fig 10.16 in Prince



VGG from 2014

Fig 10.17 in Prince

Data Augmentations

Basic:

```
146 augmentations = [
147     autocontrast, equalize, posterize, rotate, solarize, shear_x, shear_y,
148     translate_x, translate_y
149 ]
150
151 augmentations_all = [
152     autocontrast, equalize, posterize, rotate, solarize, shear_x, shear_y,
153     translate_x, translate_y, color, contrast, brightness, sharpness
154 ]
155
```

Code in PixMix Repo

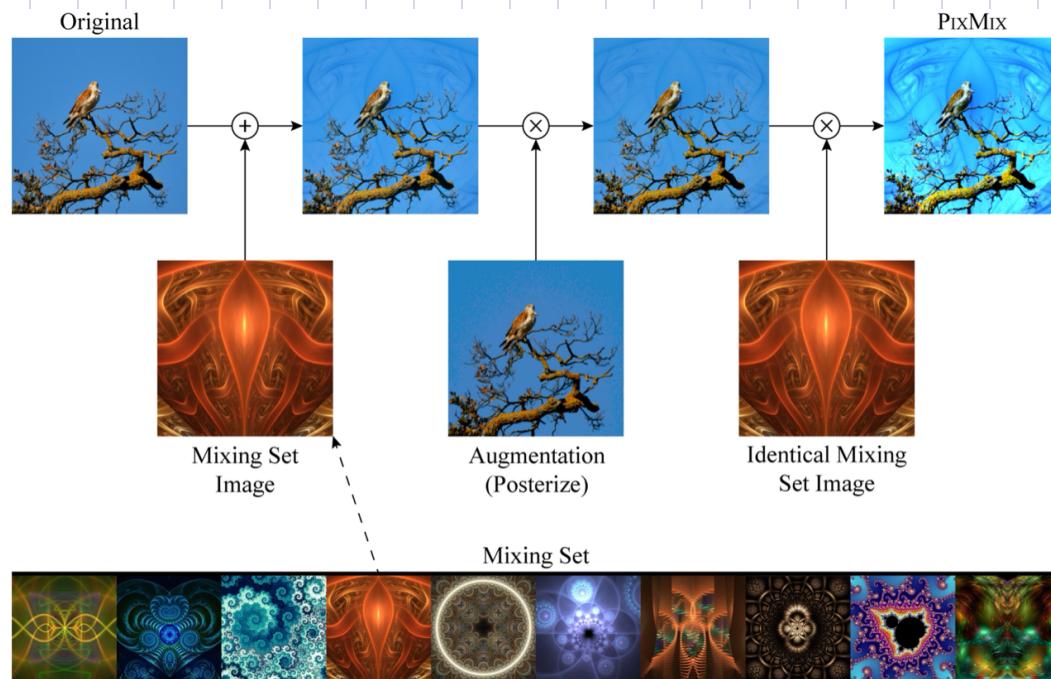
Hendrycks, et. al. CVPR
2022

More aggressive: gain robustness



PixMix Paper

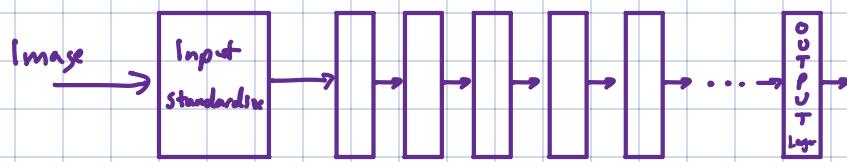
Psychedelic



PixMix Paper

Training: Stability and Effectiveness

Normalization Layers

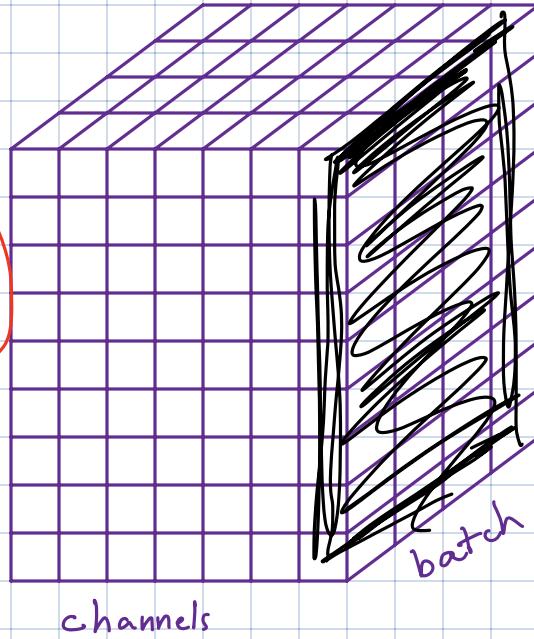


Idea: RMS Norm Layer

$$\text{d-dim } \left\{ \vec{h}_e \right\} \xrightarrow{\text{RMS Norm}} \tilde{h}_e \quad \text{where } \|\tilde{h}_e\|_{\text{RMS}} = 1$$

$$\tilde{h}_e = \frac{\vec{h}_e}{s} \quad \text{where } s = \sqrt{\frac{1}{d} \sum_i h_e^{(i)} + \epsilon}$$

$$= \frac{\vec{h}_e}{\|\vec{h}_e\| + \epsilon}$$



Batch Norm

Average over space & batch

Let \mathcal{B} be what is averaged over

$$m = \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} h_{in}$$

$$s = \sqrt{\frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} (h_{in} - m)^2}$$

$$h_{out}^{(i,j)} = \gamma \left(\frac{h_{in}^{(i,j)} - m}{s + \epsilon} \right) + \delta$$

↑
[constant]

