# Last block:

Build an encoder and a decoder to:

Accept noisy image, produce clean version

By:

Constructing loss
Applying SGD to get minimal loss on training data

# Housekeeping

There is a great deal of housekeeping:

If you build your own from scratch, you will suffer (but learn) and it won't work all that well (missing tricks)

Use an API (list in notes)

Build from someone else's code (Github, Huggingface)

# Missing tricks:

Scales and color

Blurred outputs and skip connections

Bad gradients and residual connections

Optimization tricks

# Scale and color

### Input scales:

Helpful to scale, offset pixel values so that mean is zero and standard deviation is one

# Output scales:

Real pixels are in range [0, 1]

How do we force output to be like this?

## Options:

map value to range

penalize over/under

bit of both

# Value to range

#### Terrible idea:

You might think that using

$$f(x) = \text{ReLU}(x) - \text{ReLU}(x-1)$$

would be a good idea (because it maps any x to the range [0,1]). In fact, it is a terrible idea, because once the output is outside that range, there is no gradient to push it back into the range (Section 16.3.1; **exercises** ).

# NO GRADIENT when you need it!

# Value to range

## Sigmoid:

$$sigmoid(x) = \frac{e^x}{1 + e^x}$$

but this creates gradient problems when output close to zero or one

### Tanh:

x -> [-1, 1] and easy to get to [0, 1] similar gradient problems

# **Penalties**

A second strategy is to accept the output of the convolutional layer, but penalize values that appear outside the range you want with a loss term. For example, apply the loss

$$L_{ud}(x) = x^2 \mathbb{I}_{[x<0]} + (x-1)^2 \mathbb{I}_{[x>1]}.$$

## Combination

It is often useful to mix these strategies. For example, a final layer that applies

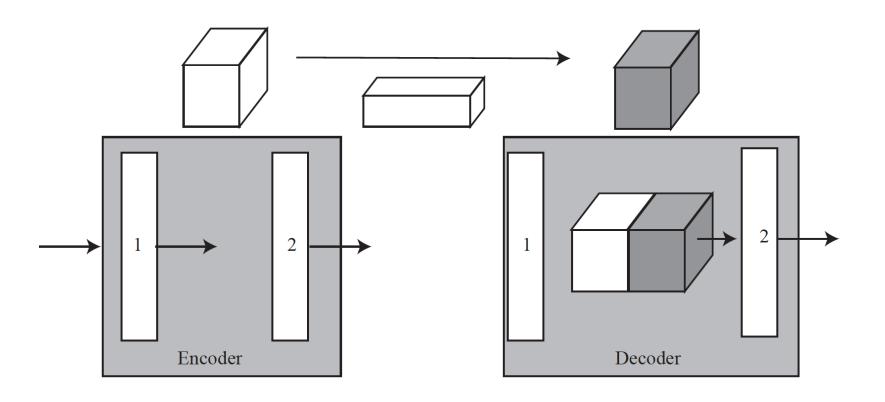
$$f(x) = a(\tanh(x) + b)$$

will map x to the range (a(b-1), a(b+1)). If you choose b < 1 and a > 1/(b+1), then the output can be below zero (but not much) and above one (but not much). Further, the gradients won't be too small at zero or one. You can then push the outputs to be in that range with a penalty term.

# Skip Connections:

An important difficulty presented by stacking many convolutional layers is that any feature produced by the encoder necessarily depends on a fairly large receptive field. This can make it difficult to produce reconstructions with sharp edges. A feature that depends on a very small neighborhood could provide enough information to place an edge accurately – for example, report the gradient of the image. If the receptive field is large, constructing a very local feature that isn't somewhat smoothed will require a set of weights that ignores many or most of the pixel values in the receptive field, which will be difficult to achieve. However, features with large receptive fields may be necessar to denoise, because they can observe long-range trends in the image.

# Skip Connections:



## Normalization:

Numbers with large magnitude in a neural network cause problems. Imagine some input to some unit is big and the weight applied to that input is small. Then a single gradient step could cause the weight to change sign, and the ReLU might cause the corresponding output to swing between strongly positive and zero. This can cause training problems, because the gradient will be a poor predictor of what will actually happen to the output. Ideally, relatively few values at the input of any layer will have large absolute values. A new layer, sometimes called a batch normalization layer, can be inserted between two existing layers to ensure this happens.

# **Batch normalization**

Write  $\mathcal{I}$  for the input of this layer, which is a  $X \times Y \times F$  block of features, and  $\mathcal{O}$  for its output, which is a block of features of the same dimension. The layer has two vectors of parameters,  $\gamma$  and  $\beta$ , each of dimension F. Write  $\gamma_i$  for the i'th component of  $\gamma$ , etc. Assume we know the mean  $(m_k)$  and standard deviation  $(s_k)$  of each feature in  $\mathcal{I}$  computed over the whole dataset and over the spatial dimensions. Write  $\epsilon$  for a small positive number chosen to avoid divide-by-zero. The data block  $\mathcal{U}$ , with ijk'th component

$$\mathcal{U}_{ijk} = \frac{(\mathcal{I}_{ijk} - m_k)}{(s_k + \epsilon)}$$

will tend to have small magnitude numbers in it, both positive and negative. The mean of each feature in this block should be about zero, because it is close to the mean over all blocks. The standard deviation of each feature in this block should be about one, because it is close to the standard deviation over all blocks. Now compute

$$\mathcal{O}_{ijk} = \gamma_k \mathcal{U}_{ijk} + \beta_k$$

and notice that  $\mathcal{O}$  could be the same as  $\mathcal{I}$  (set  $\gamma_k = s_k$  and  $\beta_k = m_k$ ). The output of this layer is a differentiable function of  $\gamma$  and  $\beta$ , which can be adjusted to achieve the best performance.

### **Batch normalization**

But we don't know the mean and standard deviation! And it keeps changing during training anyhow!

#### Strategy:

Use the mean and standard deviation of the last batch at training time At test time, use constants (last seen mean/sd)
API looks after this housekeeping

#### Landmine:

You have to tell the networks when you switch from train to test If you don't they work badly