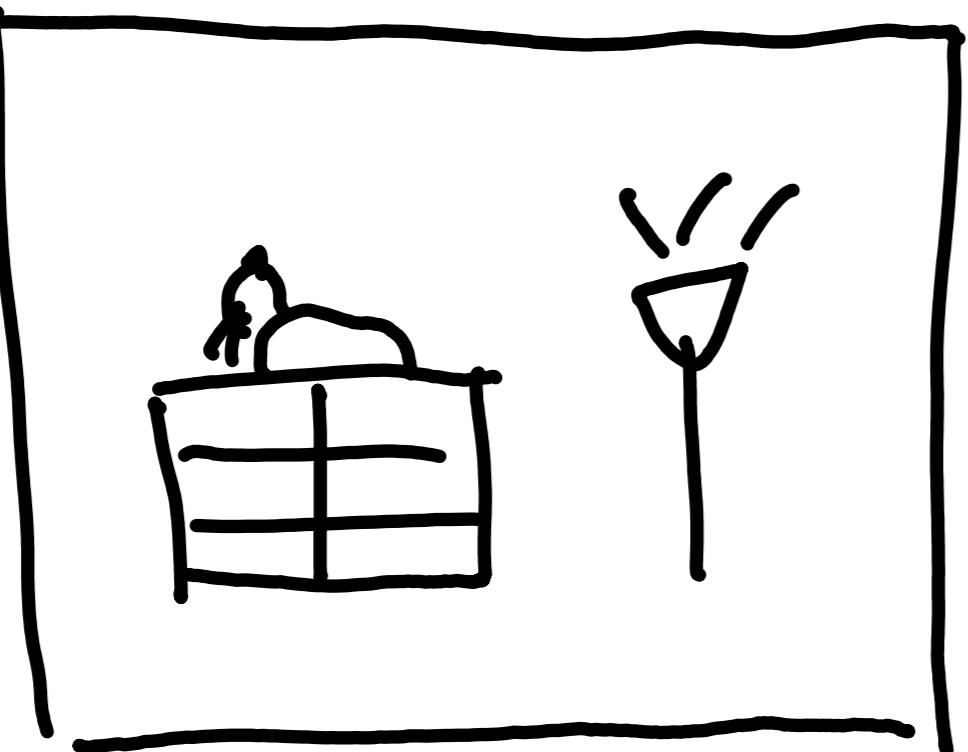


Convolutional Neural Networks



Does this image
contain a cat?

1. Translation invariance
(respond similarly everywhere)
2. Locality
(only consider a "small" region)

$$h = u + w_x$$

$$h_m = u_m + \sum_n w_{m,n} x_n$$

$$H_{i,j} = U_{i,j} + \sum_k \sum_\ell w_{i,j,k,\ell} x_{k,\ell}$$

↪
 $i,j = m, k,\ell = n$
 ↪

$$H_{i,j} = U_{i,j} + \sum_a \sum_b V_{i,j,a,b} X_{i+a, j+b}$$

↪ can be
 positive or negative
 to go over whole image

Translation invariance:

$$H_{i,j} = U + \sum_a \sum_b V_{a,b} X_{i+a, j+b}$$

Locality:

$$H_{i,j} = U + \sum_{a=-\Delta}^{\Delta} \sum_{b=-\Delta}^{\Delta} V_{a,b} X_{i+a, j+b}$$

"kernel"
"filter"

How much is
U present
at i, j?

Convolution!

Channels: Images are 2D grids of pixels with one or more channels (R G B)

In convnets, feature representations (activations) have channels.

$$H_{i,j,d} = U_{c,d}^{??} + \sum_{a=-\Delta}^{\Delta} \sum_{b=-\Delta}^{\Delta} \sum_c V_{a,b,c,d} X_{i+a, j+b, c}$$

↑
output channel

↑
sum over input channels

↑
input channel

$$\begin{array}{|c|c|c|} \hline 0 & 1 & 2 \\ \hline 3 & 4 & 5 \\ \hline 6 & 7 & 8 \\ \hline \end{array} *_{\text{Conv}} \begin{array}{|c|c|} \hline 0 & 1 \\ \hline 2 & 3 \\ \hline \end{array} = \begin{array}{|c|c|} \hline 19 & 25 \\ \hline 37 & 43 \\ \hline \end{array}$$

X V H
 input filter output
 (one channel) or kernel

Notes:

1. Output is smaller than the input.

$$\text{output size} = (\frac{n_h - k_h + 1}{\text{input height}}) \times (\frac{n_w - k_w + 1}{\text{input width}})$$

$$= (\frac{n_h - k_h + 1}{\text{kernel height}}) \times (\frac{n_w - k_w + 1}{\text{kernel width}})$$

2. "Receptive field": Region in the input that a given position in the output depends on

Padding: Add zeros around the border of the input so that the output is the same size.

Padded input:

$$\begin{matrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 2 & 0 \\ 0 & 3 & 4 & 5 & 0 \\ 0 & 6 & 7 & 8 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{matrix}$$

*

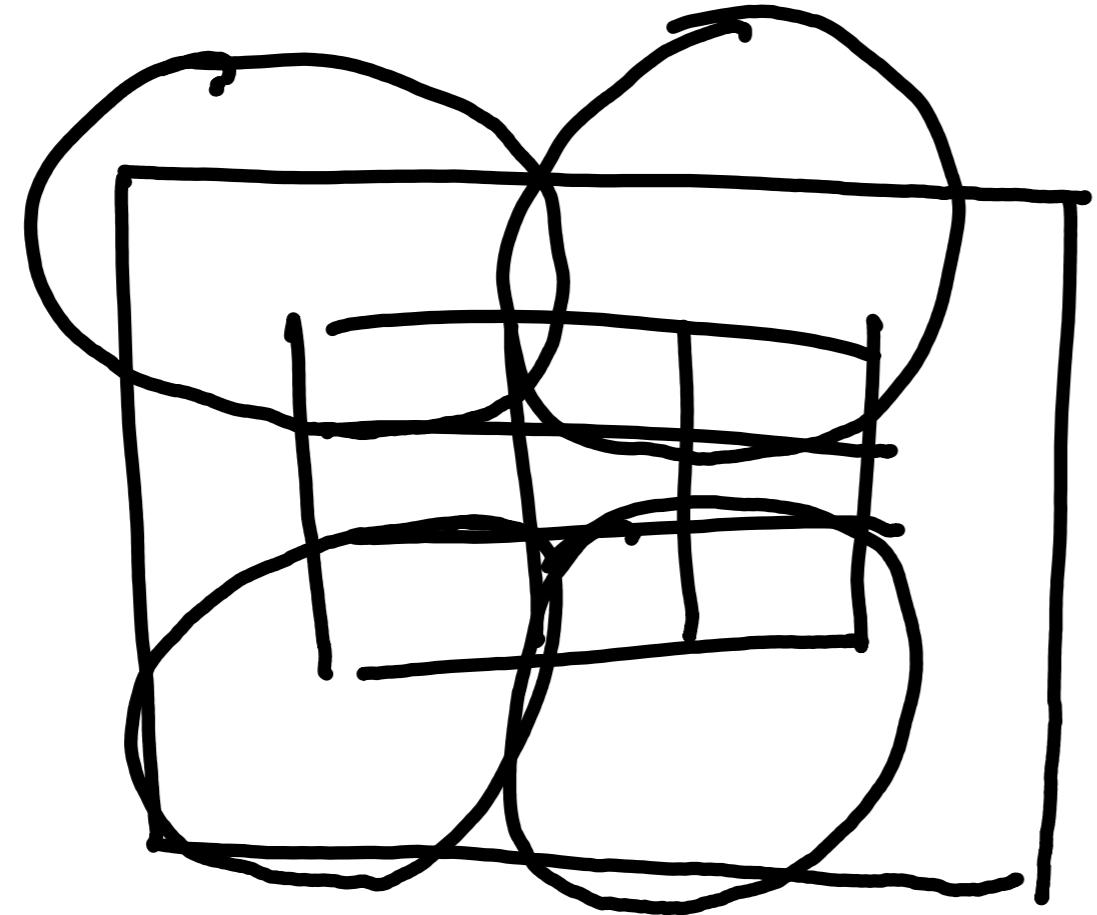
$$\begin{matrix} * & * & * \\ * & * & * \\ * & * & * \end{matrix}$$
 =

$$\begin{matrix} * & * & * \\ * & * & * \\ * & * & * \end{matrix}$$

$k \times k$ kernel
 k odd
 same-size

total padding on axis : $k-1$
on each edge, we need $\frac{k-1}{2}$ zeros.

Stride:



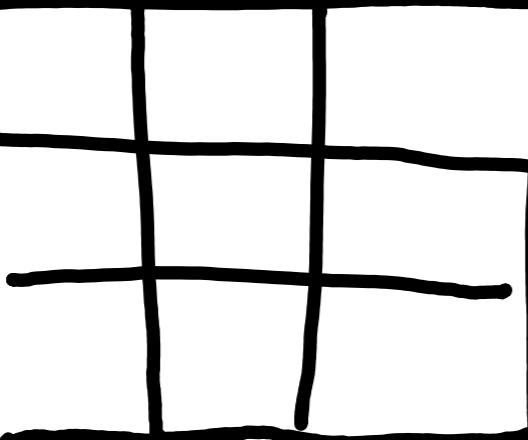
$$* \begin{array}{|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} = \begin{array}{|c|c|} \hline & \\ \hline & \\ \hline \end{array}$$

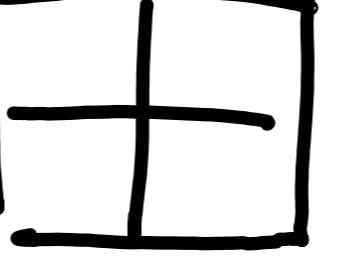
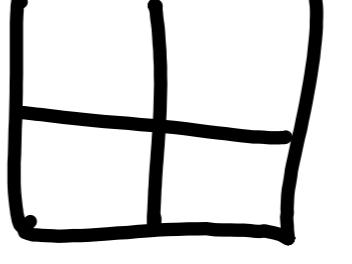
Stride of 2×3

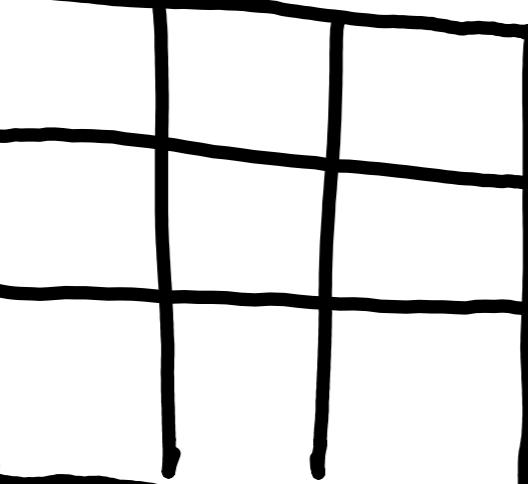
↗ amount we move horizontally
↖ amount we move vertically,

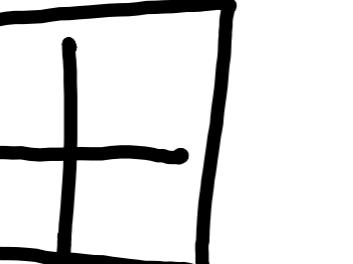
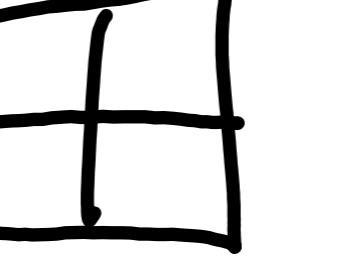
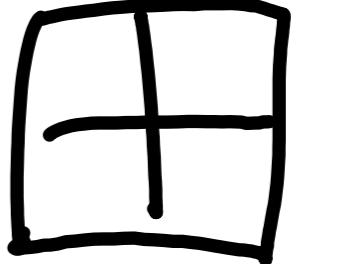
- output gets smaller
"down sampling"
- cheaper!

Channels

$$x_{\cdot, \cdot, 1}$$


$$\ast$$
 = 

$$x_{\cdot, \cdot, 2}$$


$$\ast$$
 =  + = 

- Separate kernel for each input/output channel combination
- Sum the result over input channels

Pooling: Aggregate spatial information

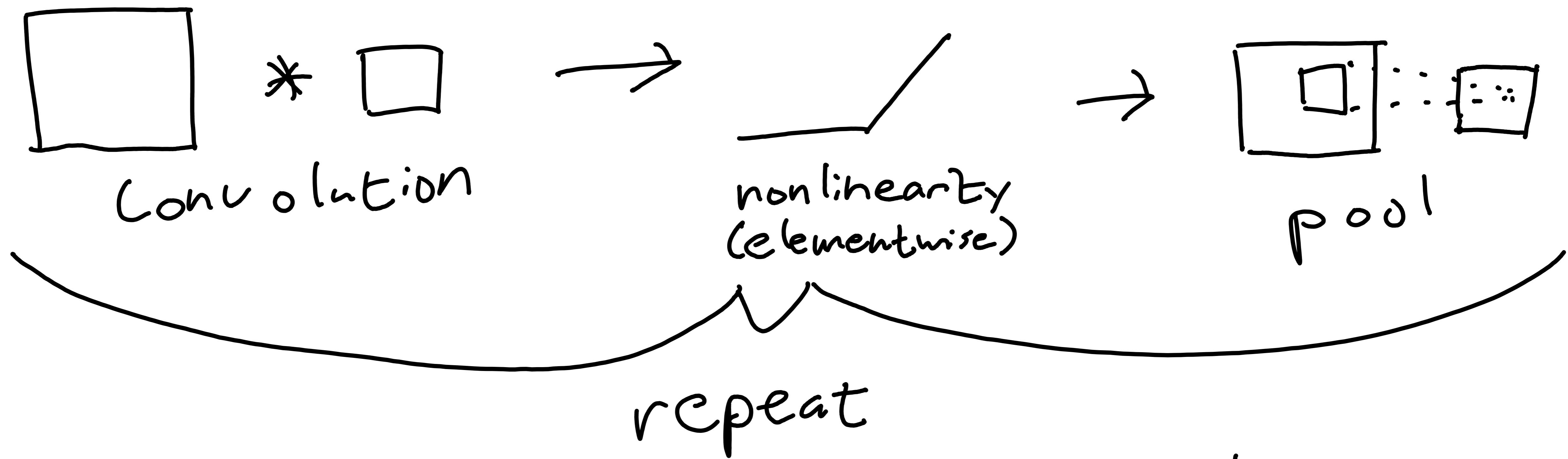
max pooling:

$$\begin{array}{|c|c|c|} \hline 0 & 1 & 2 \\ \hline 3 & 4 & 5 \\ \hline 6 & 7 & 8 \\ \hline \end{array} \rightarrow \begin{matrix} 2 \times 2 \\ \text{max pool} \end{matrix} \rightarrow \begin{array}{|c|c|} \hline 4 & 5 \\ \hline 7 & 8 \\ \hline \end{array} \quad (\text{stride } 1)$$

When performing pooling, almost always stride by pooling size

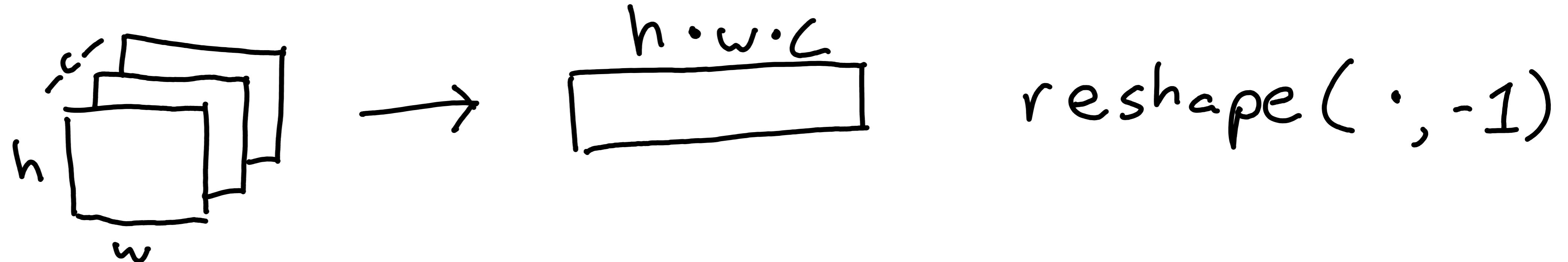
Average pooling is also not uncommon.

Simple ConvNet structure:



Eventually we need to perform classification.

1. Flatten ($7 \times 7 \times 100 \rightarrow 4900$)



2. Apply dense layers (like an MLP)

$$3. BN(x) = \frac{\gamma x - \mu_B}{\sigma_B + \epsilon} + \beta$$

*learned
params*

$$\gamma, \beta \in \mathbb{R}^{\dim(x)}$$

4. At inference time, use EMA of μ_B and σ_B computed over training.

Batch norm for images:

$$x \in \mathbb{R}^{b \times h \times w \times c}$$

$$\mu = \text{mean}(x, \text{axis}=(0, 1, 2)) \in \mathbb{R}^c$$

$$\sigma = \text{std}(x, \text{axis}=(0, 1, 2)) \in \mathbb{R}^c$$

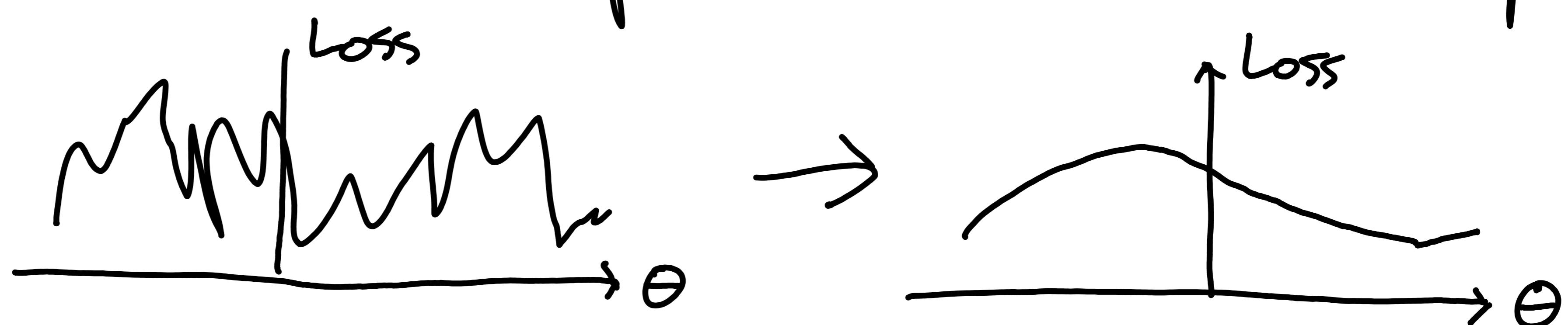
$$BN(x) = \gamma \odot \frac{x - \mu}{\sigma + \epsilon} + \beta$$

Why batch norm?

1. "Reduces internal covariate shift"
→ But in practice, this doesn't happen
2. Moves the activations toward the non-saturated region of the nonlinearity



3. Smooths the optimization landscape



4. Regularization due noisy estimate of mean and variance

Residual Connections

Say we have a model and we want to add a layer, but we only want to enlarge the set of functions it can implement.

Alternatively, we want to propagate gradients as if the layer is not there.

Say x is the output of orig. model
 $f(\cdot)$ is the new layer

"Residual connection"

$$f(x) + x$$