

MACHINE LEARNING BASICS

CMSC764 / AMSC607

LINEAR CLASSIFIERS

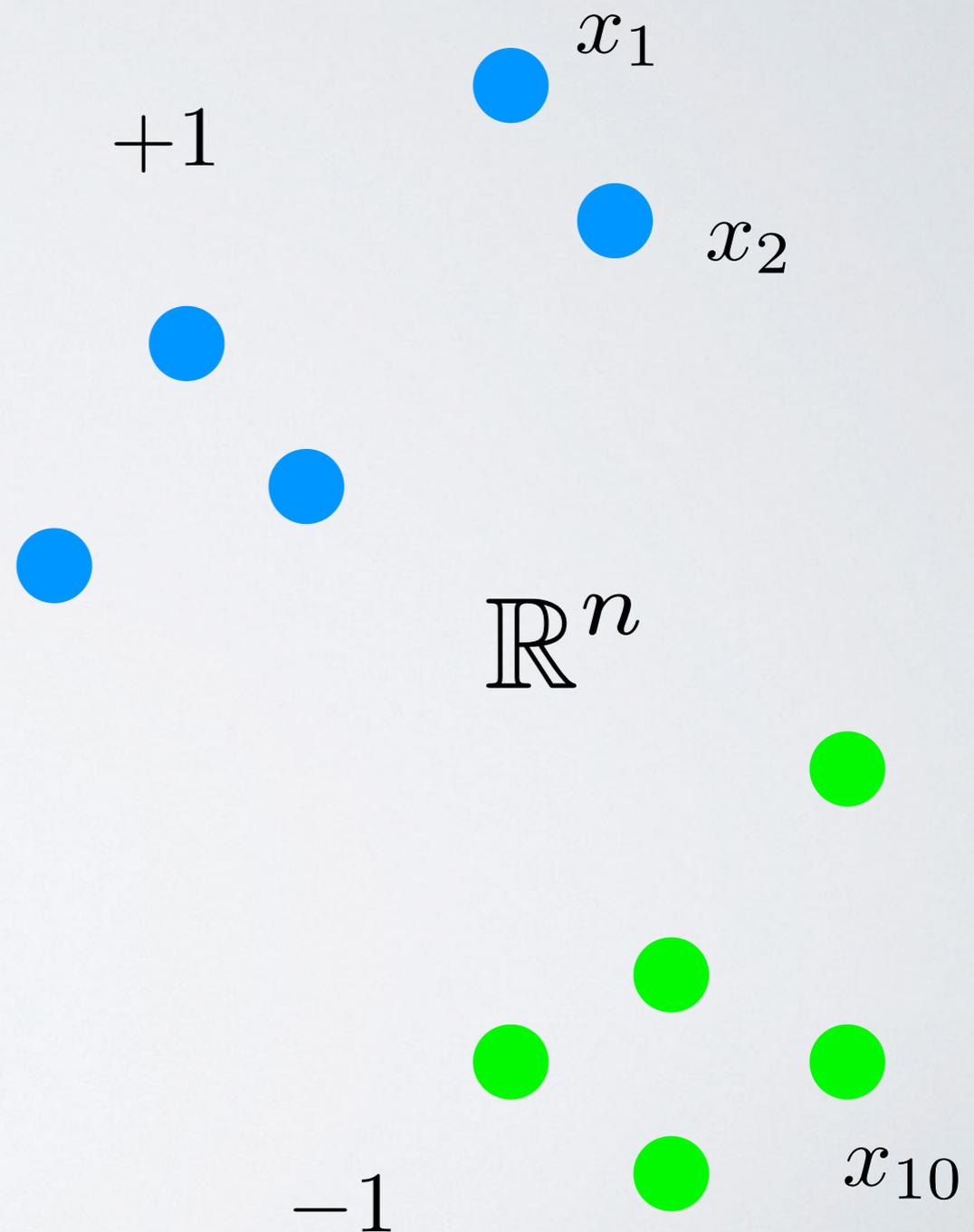
training data

feature vectors

vectors containing
descriptions of objects

labels

+1/-1 labels indicating
which “class” each vector
lies in

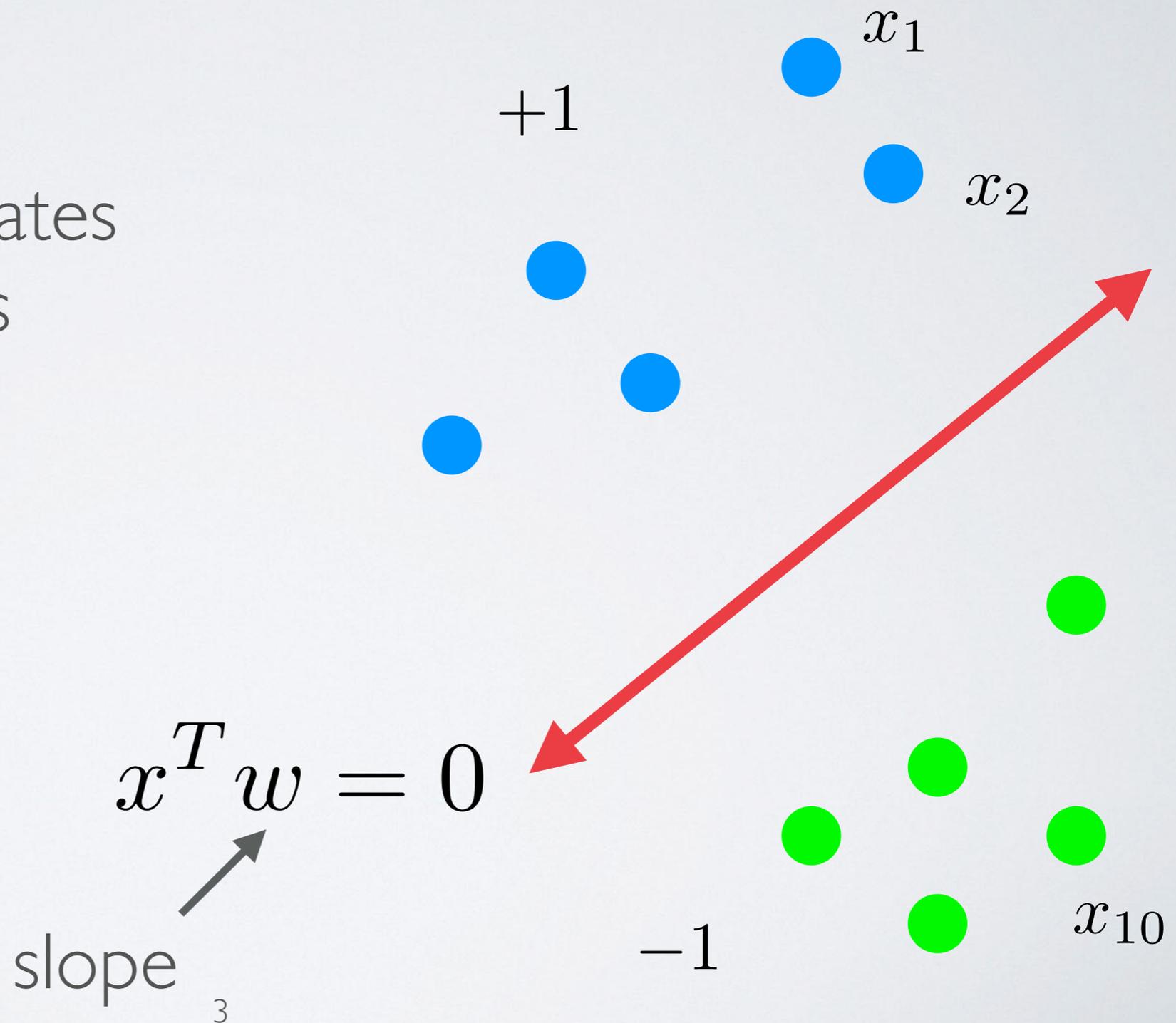


LINEAR CLASSIFIERS

goal

learn a line that separates
two object classes

what line is best?

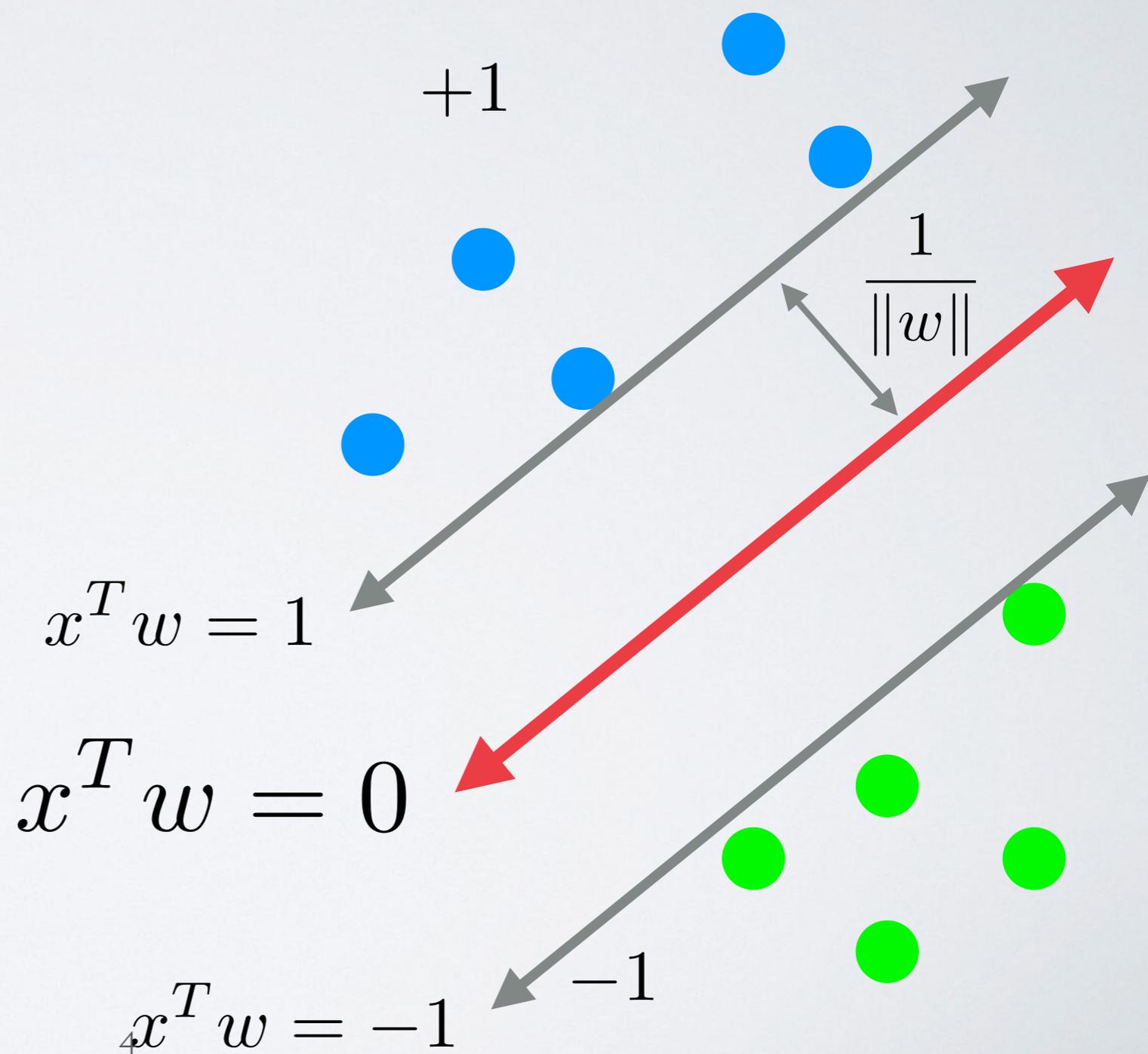


SUPPORT VECTOR MACHINE

SVM

choose line with maximum "margin"

$$\text{margin width} = \frac{1}{\|w\|}$$

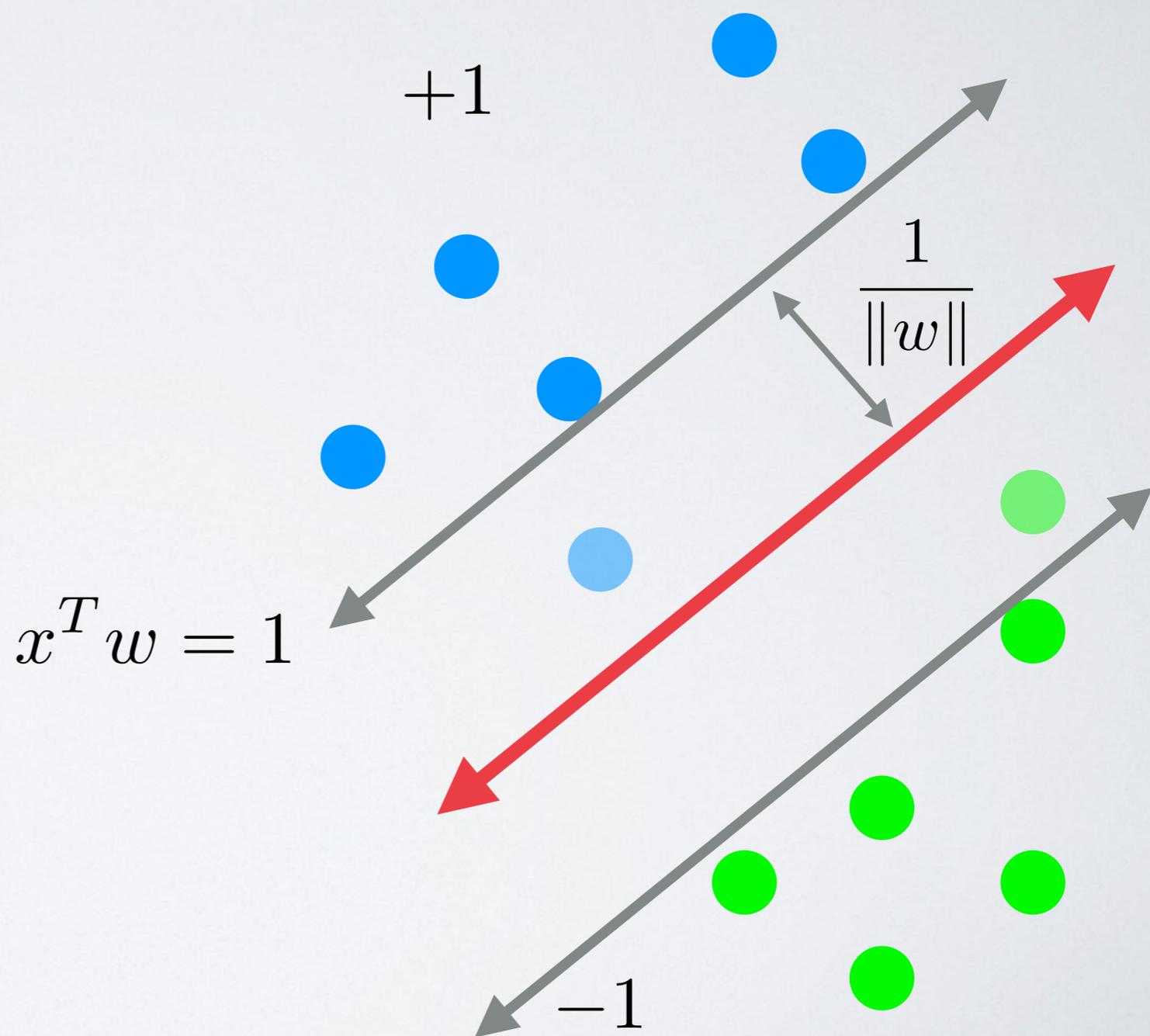
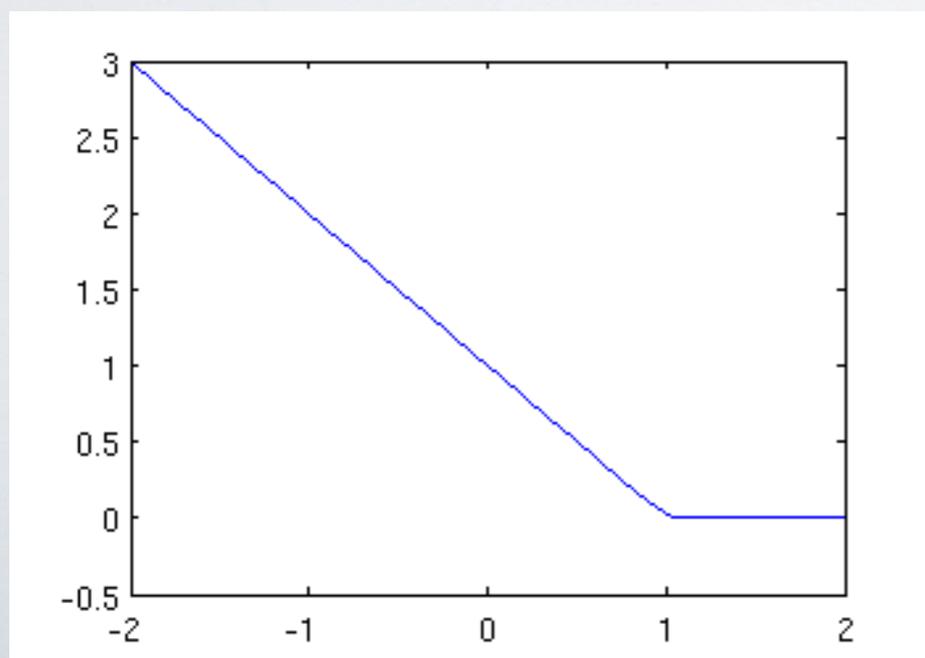


SUPPORT VECTOR MACHINE

hinge loss

penalize points that lie within the margin

$$\sum_i h(x_i^T w)$$



SUPPORT VECTOR MACHINE

hinge loss

penalize points that lie within the margin

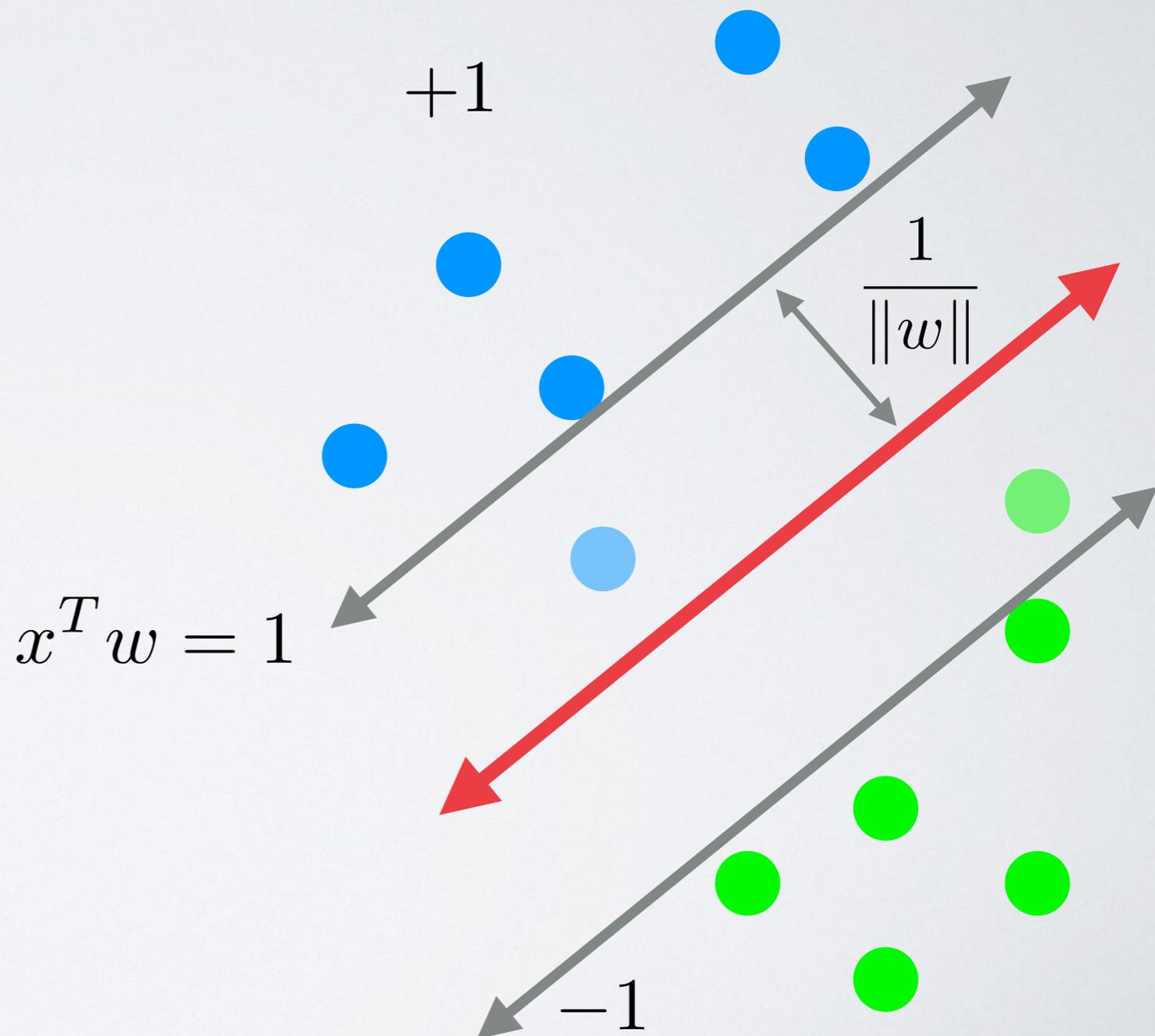
$$\sum_i h(x_i^T w)$$

For class -1 ...

$$\sum_i h(-x_i^T w)$$

In general...

$$\sum_i h(y_i \cdot x_i^T w)$$



SUPPORT VECTOR MACHINE

In general...

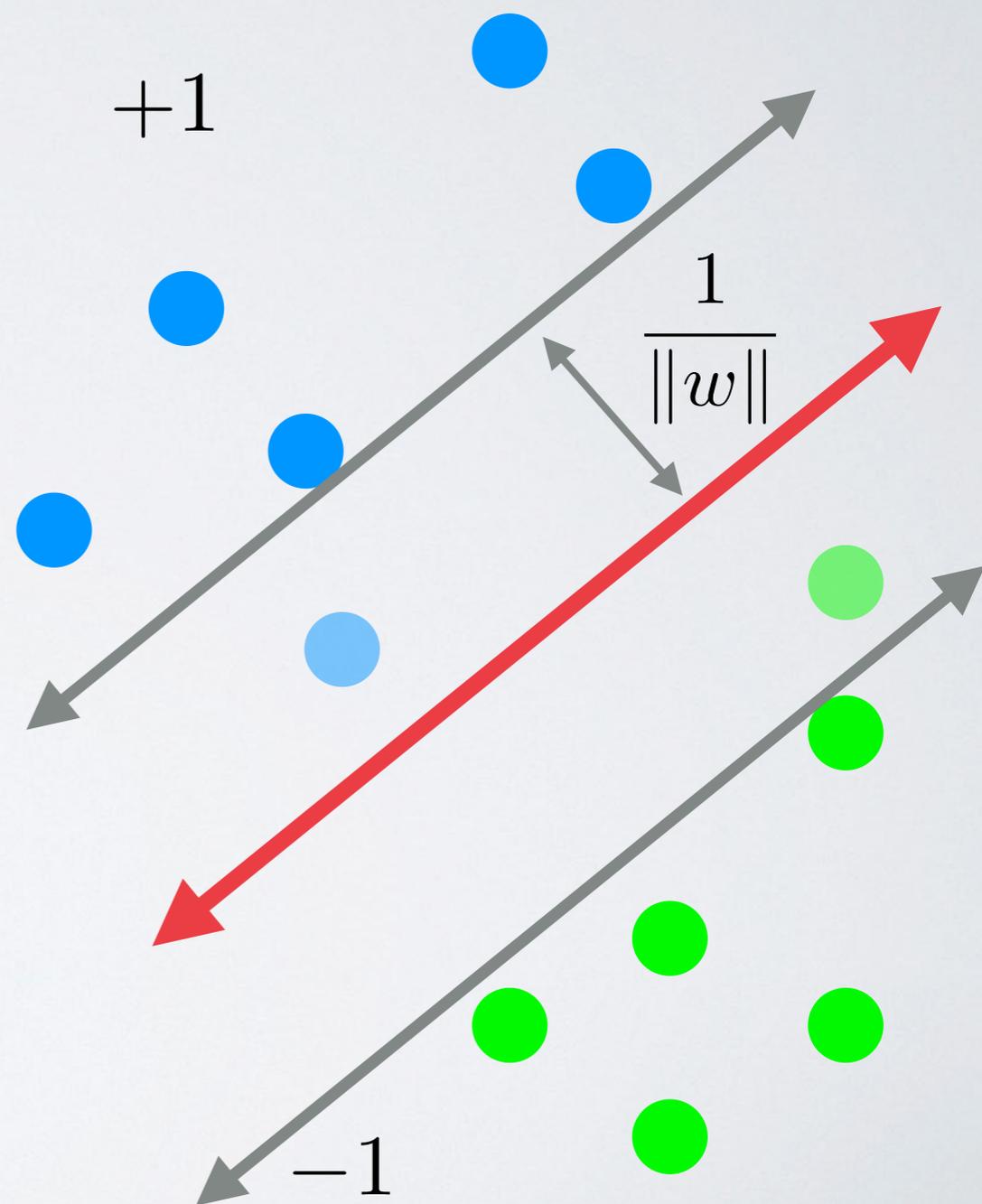
$$\sum_i h(y_i \cdot x_i^T w)$$

short-hand

$$h(\hat{X}w)$$

$$\text{margin width} = \frac{1}{\|w\|}$$

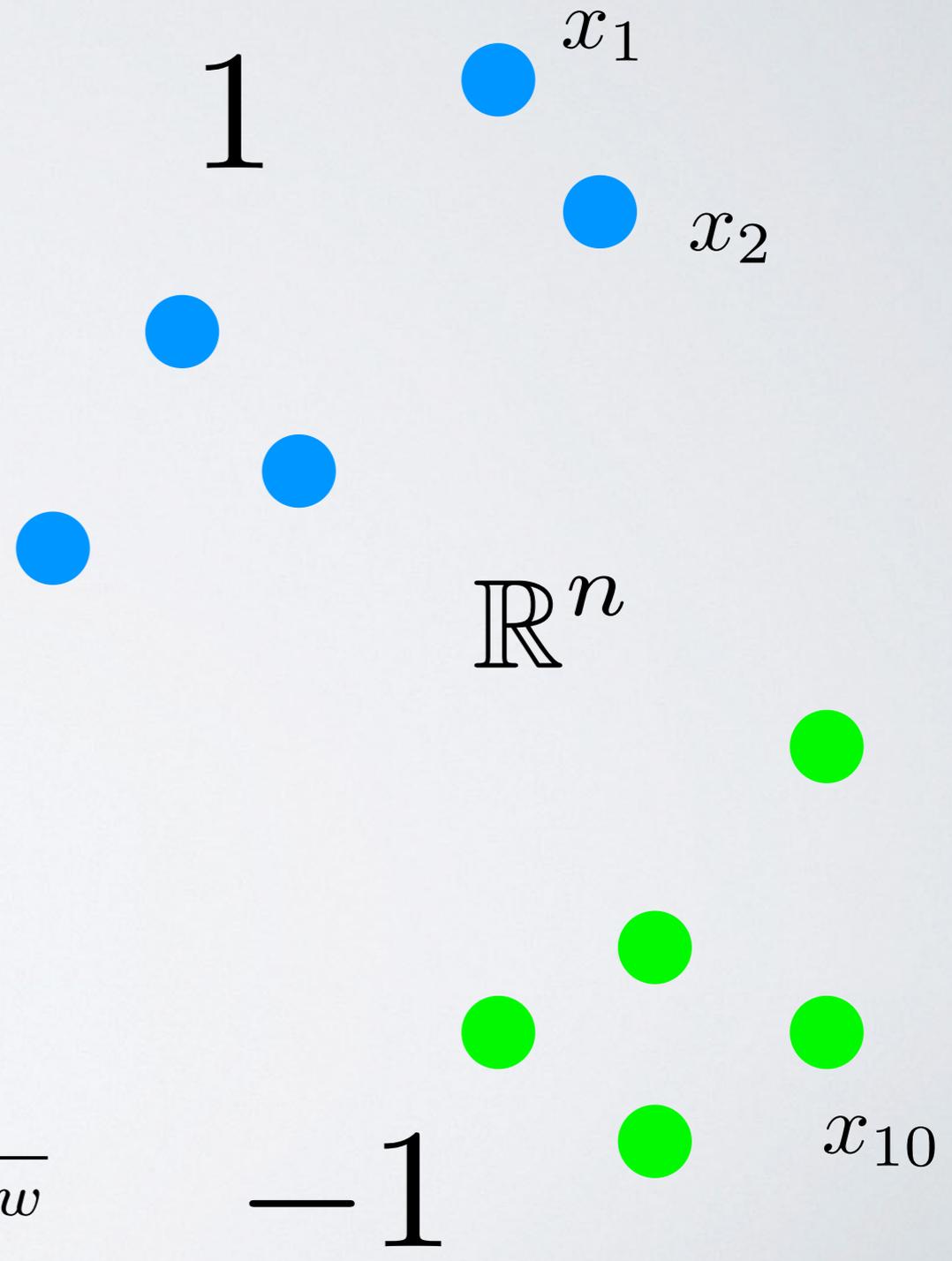
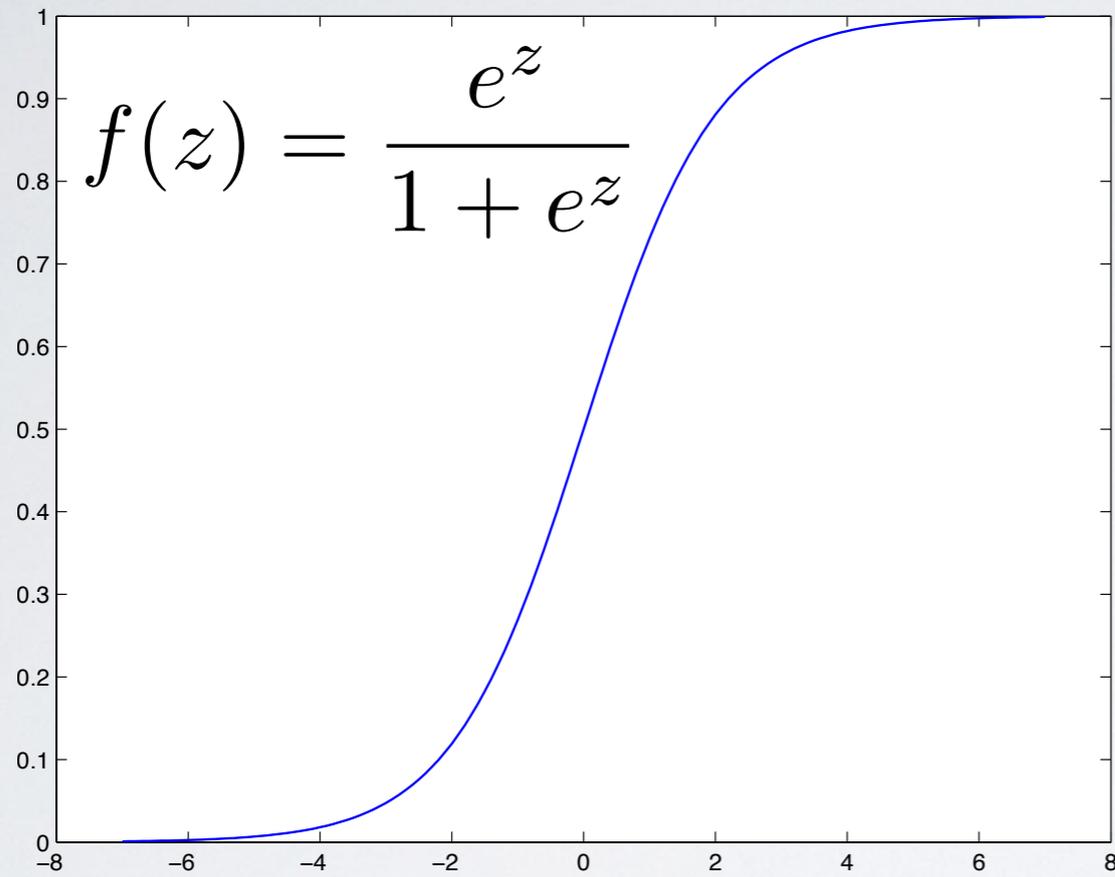
$$x^T w = 1$$



combined objective

$$\text{minimize } \frac{1}{2} \|w\|^2 + h(\hat{X}w)$$

LOGISTIC REGRESSION



$$\mathbb{P}[y = 1] = \frac{e^{x^T w}}{1 + e^{x^T w}}$$

$$\mathbb{P}[y = -1] = 1 - \frac{e^{x^T w}}{1 + e^{x^T w}} = \frac{1}{1 + e^{x^T w}}$$

LIKELIHOOD FUNCTIONS

probability of observing label
given data

case: observed $y=1$

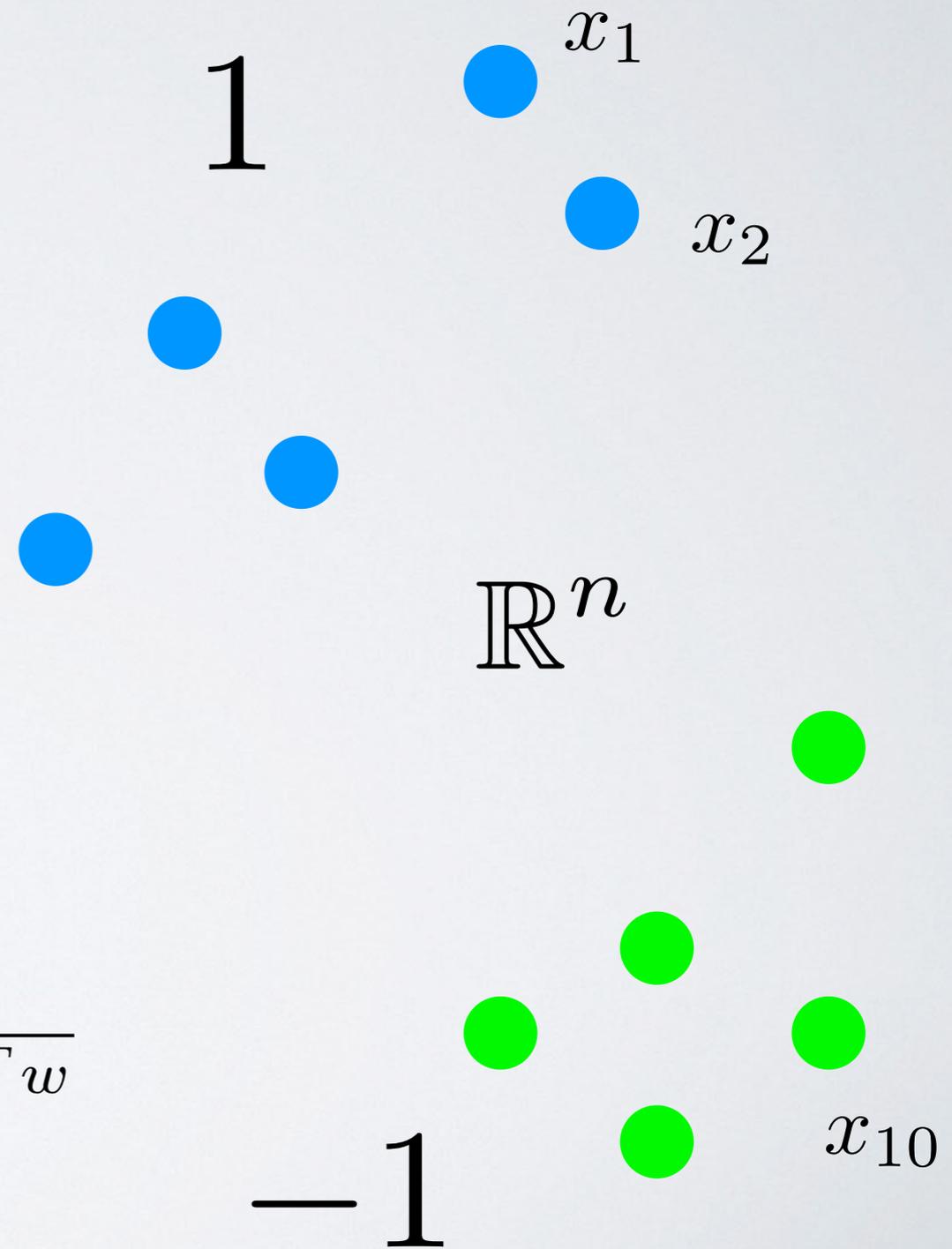
$$\mathbb{P}[y = 1] = \frac{e^{x^T w}}{1 + e^{x^T w}}$$

$$NLL(x) = \log(1 + e^{x^T w}) - x^T w$$

case: observed $y=-1$

$$\mathbb{P}[y = -1] = 1 - \frac{e^{x^T w}}{1 + e^{x^T w}} = \frac{1}{1 + e^{x^T w}}$$

$$NLL(x) = \log(1 + e^{x^T w})$$



LIKELIHOOD FUNCTIONS

probability of observing label
given data

case: observed $y=1$

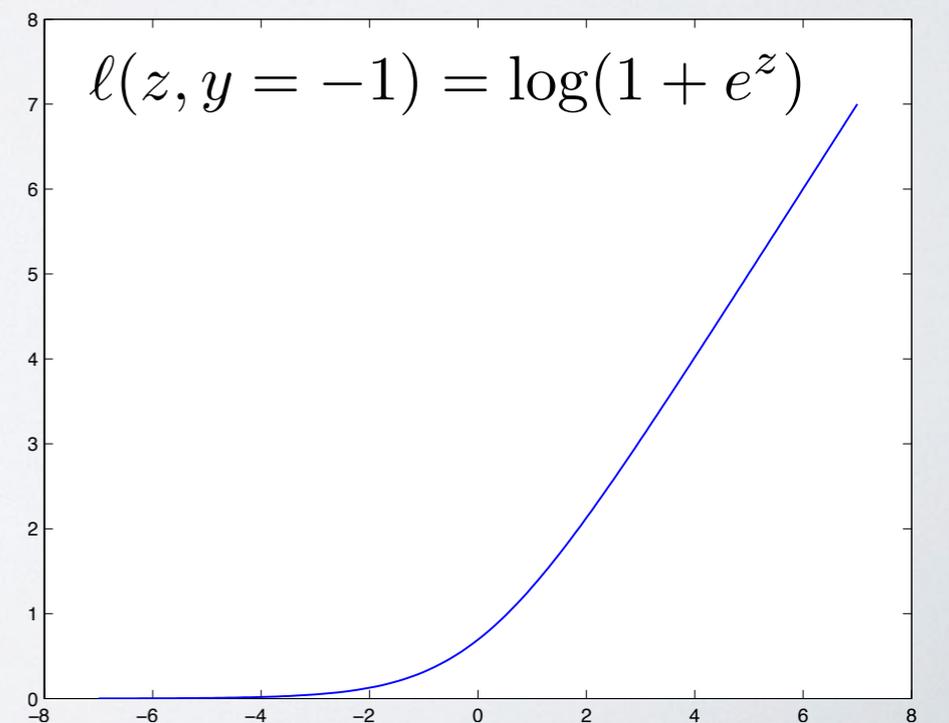
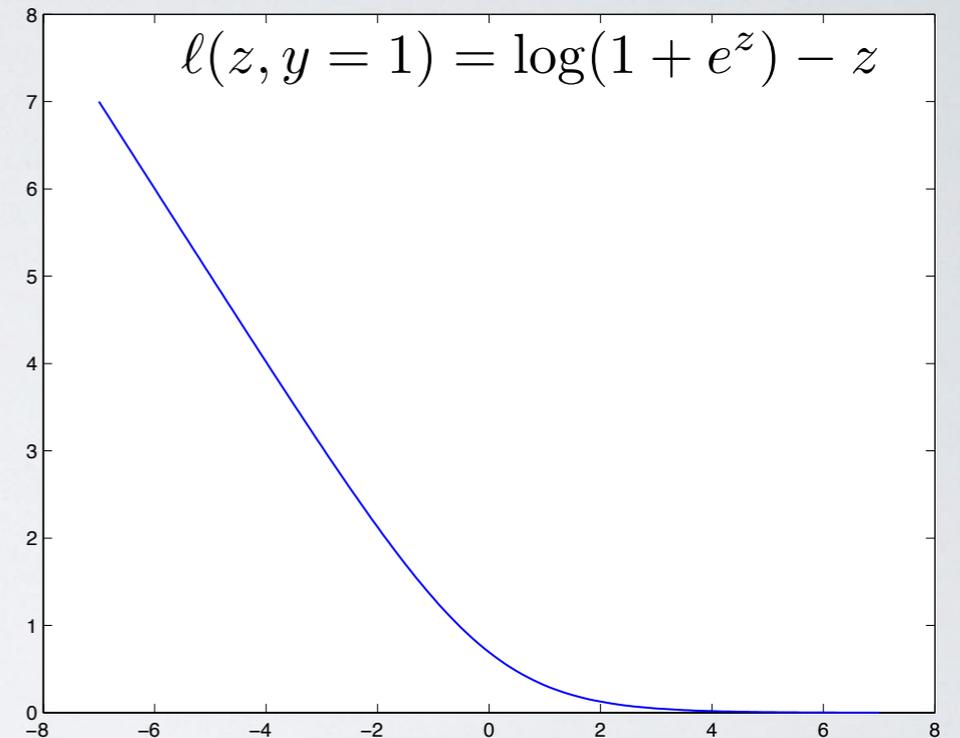
$$\mathbb{P}[y = 1] = \frac{e^{x^T w}}{1 + e^{x^T w}}$$

$$NLL(x) = \log(1 + e^{x^T w}) - x^T w$$

case: observed $y=-1$

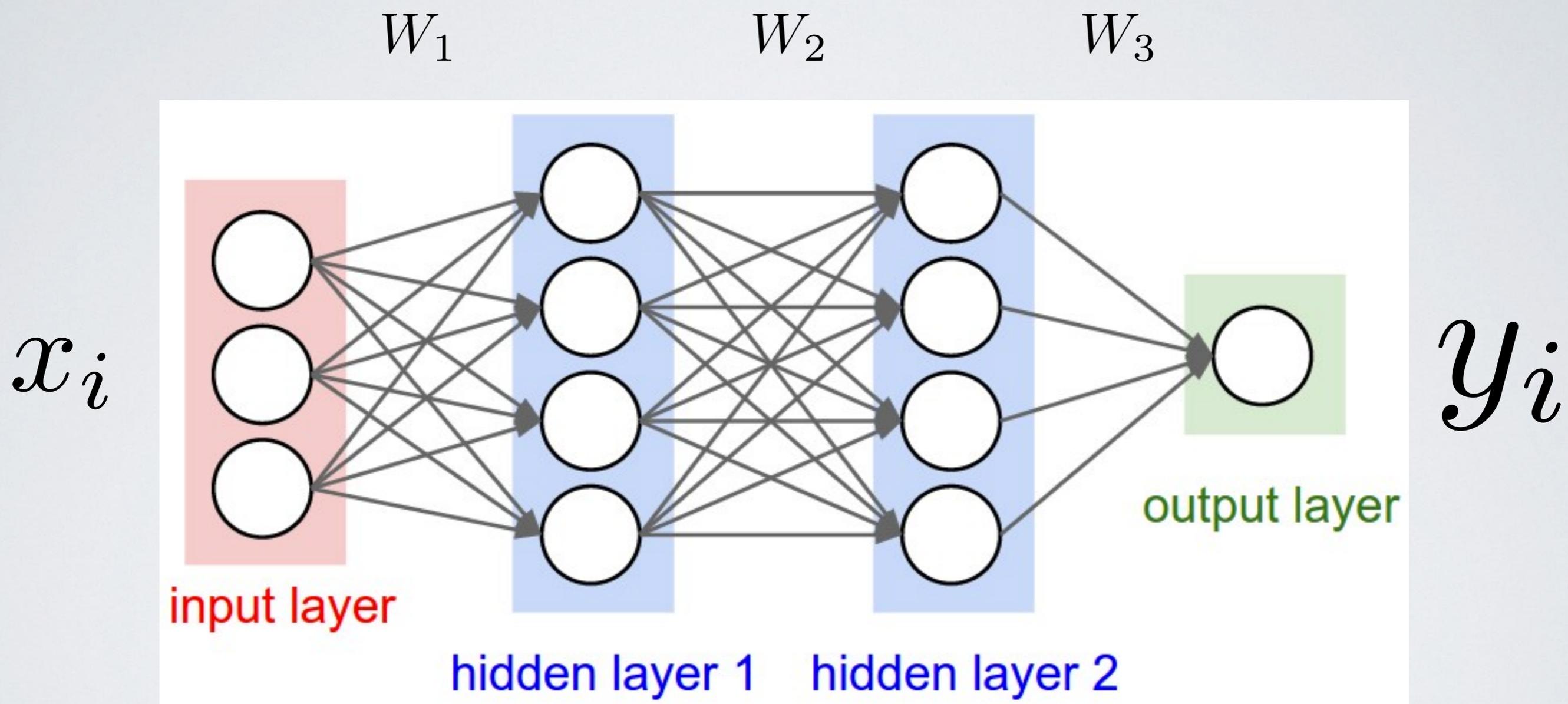
$$\mathbb{P}[y = -1] = 1 - \frac{e^{x^T w}}{1 + e^{x^T w}} = \frac{1}{1 + e^{x^T w}}$$

$$NLL(x) = \log(1 + e^{x^T w})$$

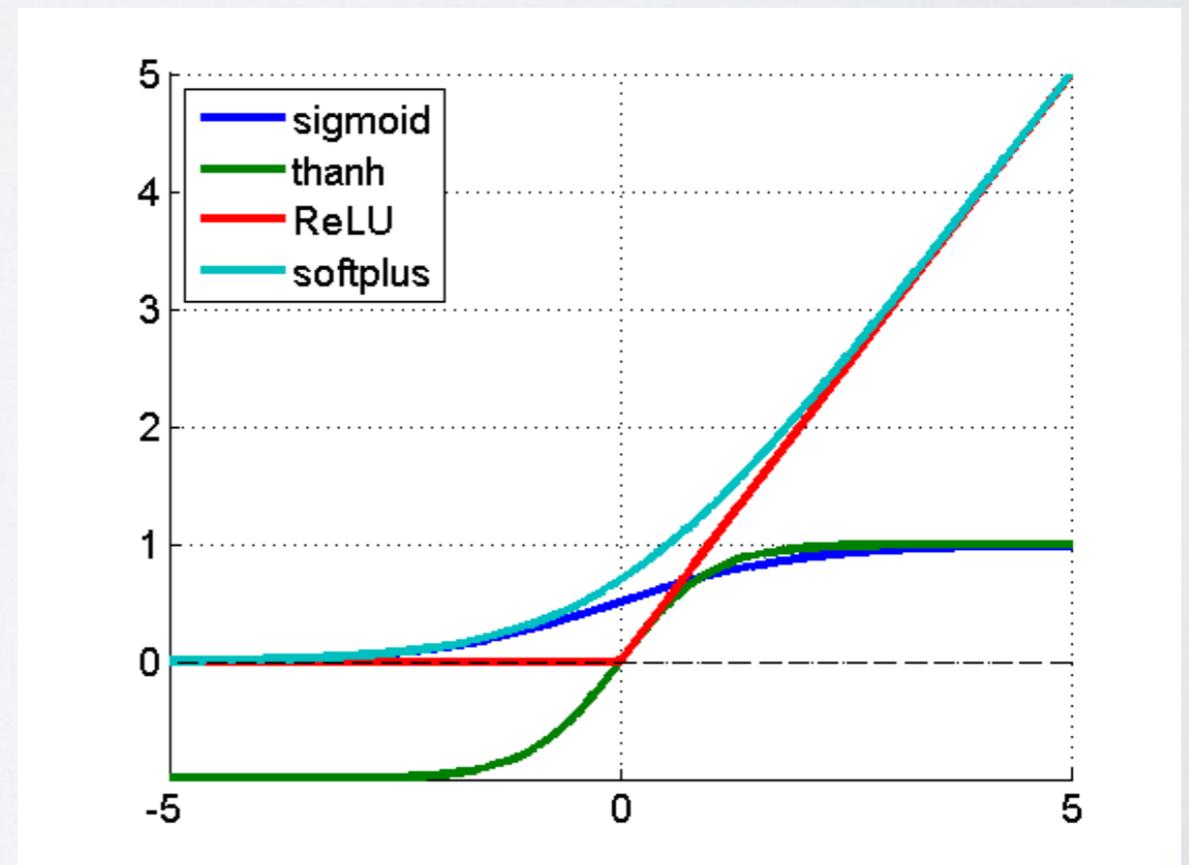
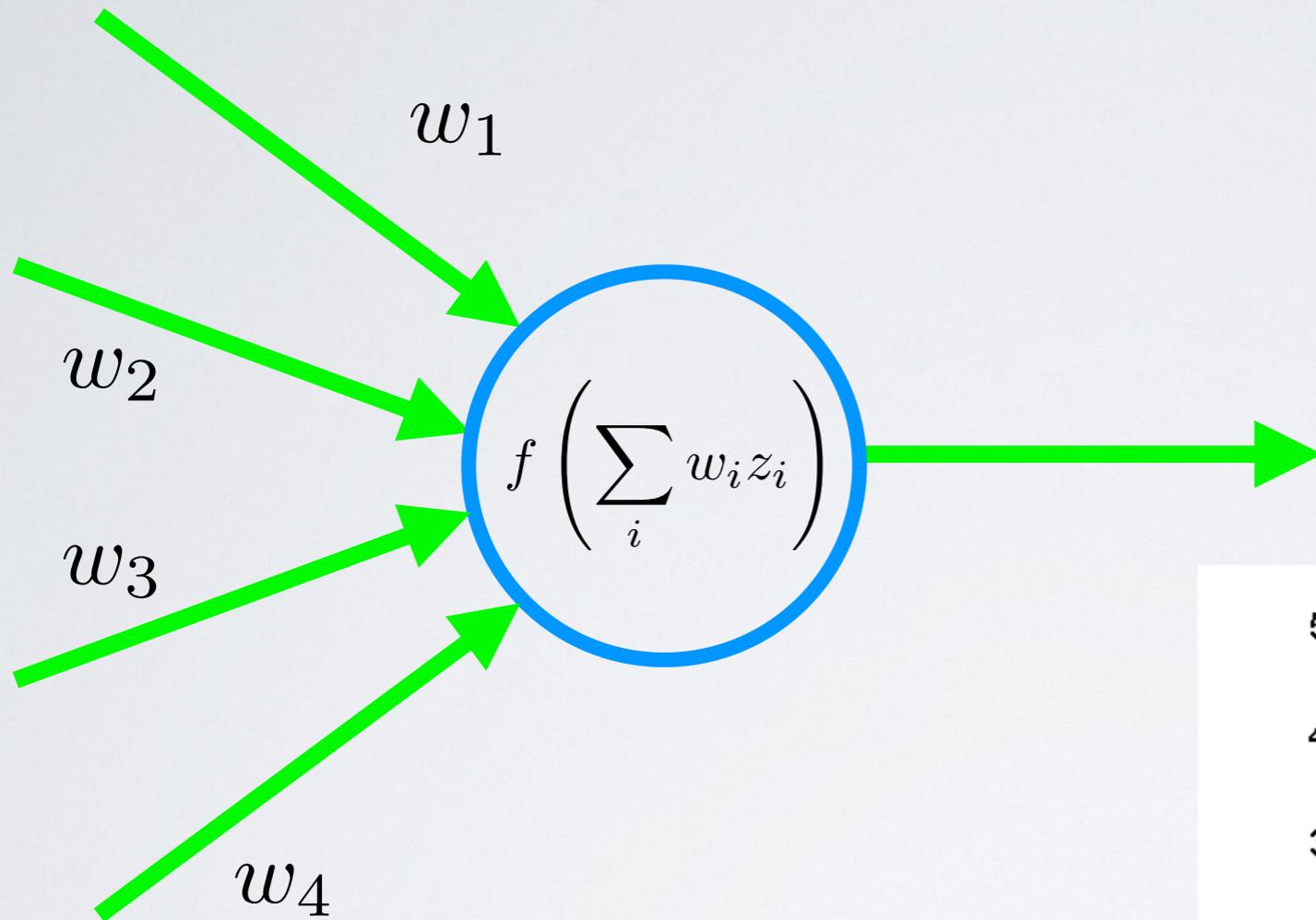


NEURAL NET BASICS

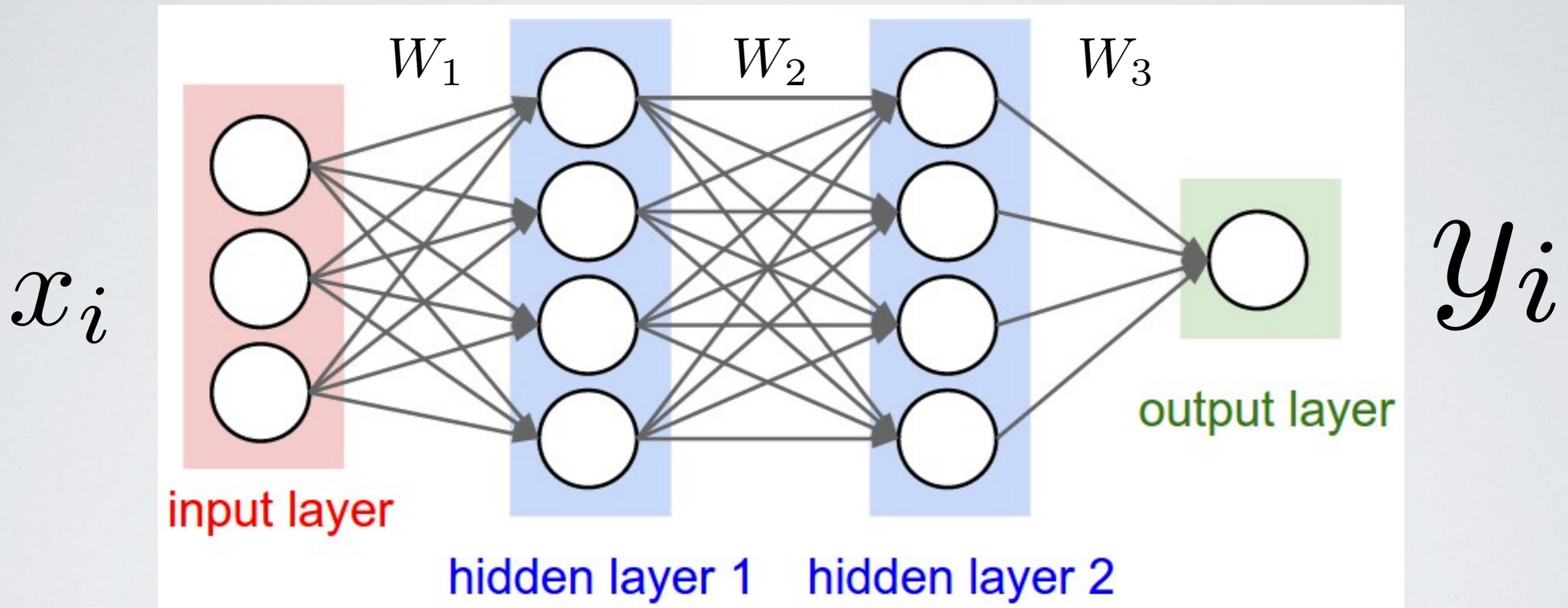
NEURAL NETWORKS



ARTIFICIAL NEURON



NEURAL NETWORKS



$$f(f(f(x_i W_1) W_2) W_3) = y_i$$

TRAINING

$$f(f(f(x_i W_1) W_2) W_3) = y_i$$

least-squares loss

$$\ell(z, y) = (z - y)^2$$

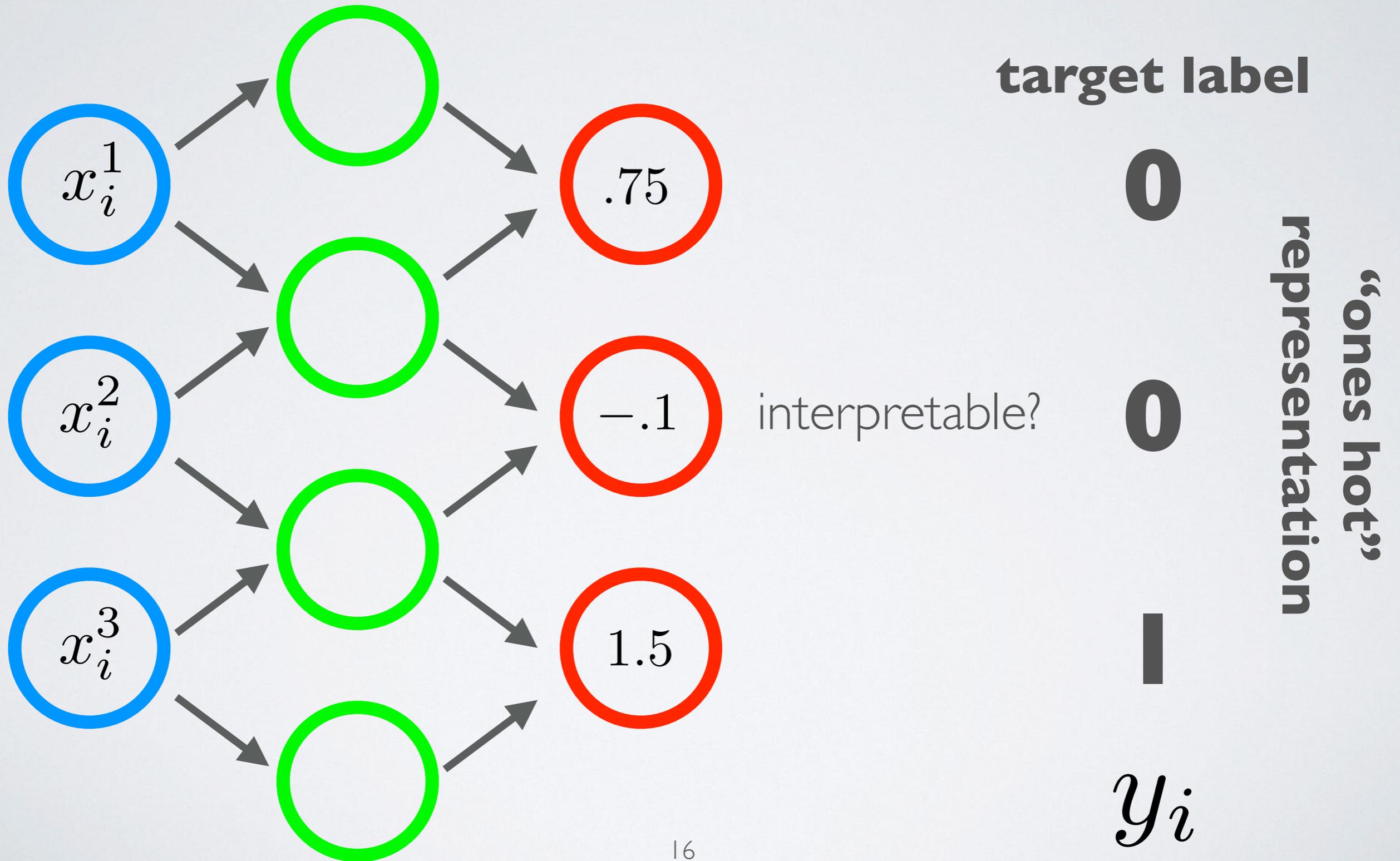
$$\min_{W_1, W_2, W_3} \sum_i \|f(f(f(x_i W_1) W_2) W_3) - y_i\|^2$$

“cross entropy”/log-likelihood loss

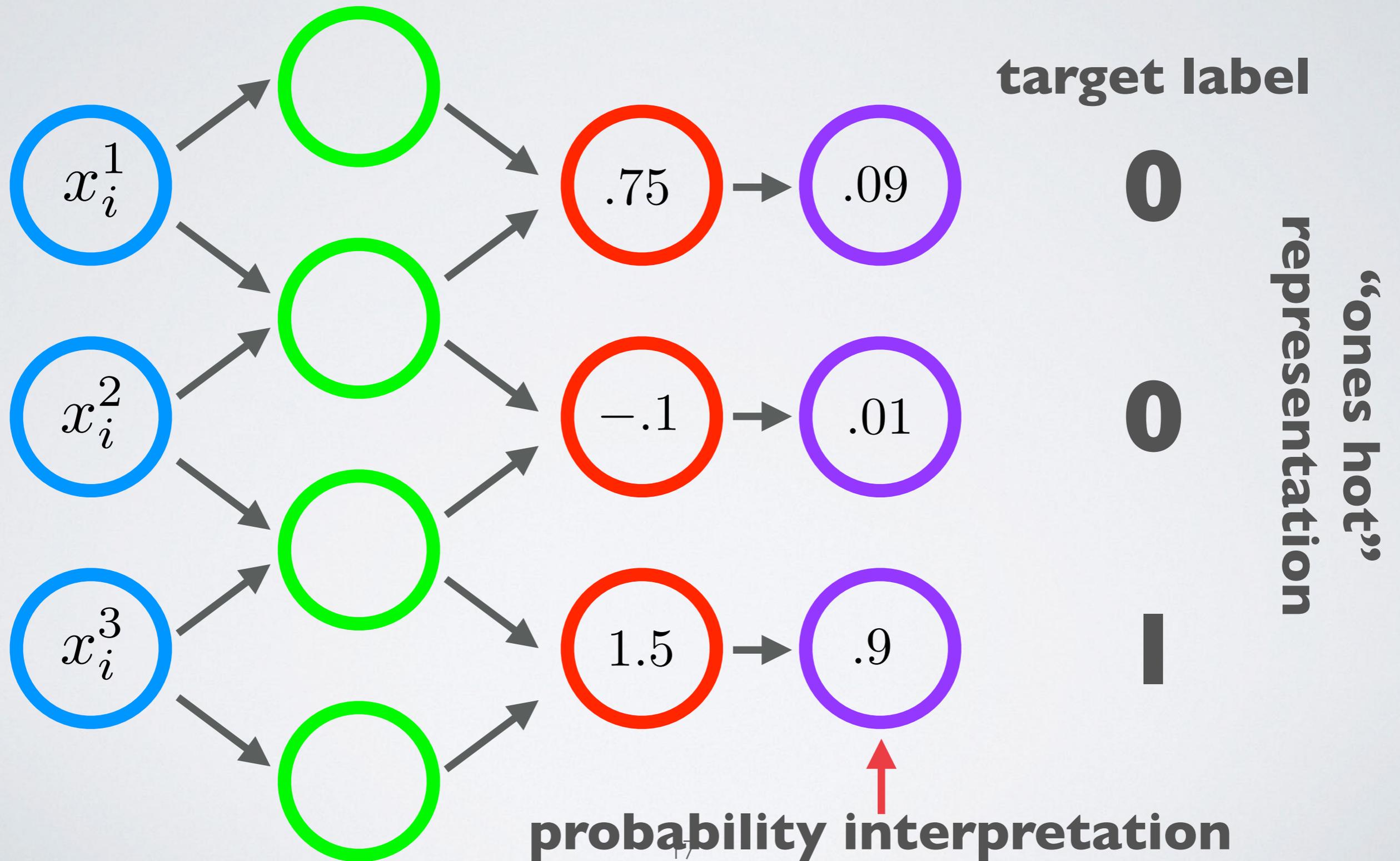
$$l(z, y) = \begin{cases} -\log(z), & y = 1 \\ -\log(1 - z), & y = -1 \end{cases}$$

$$\min_{W_1, W_2, W_3} \sum_i \ell(f(f(f(x_i W_1) W_2) W_3), y_i)$$

MULTI-CLASS NETWORK

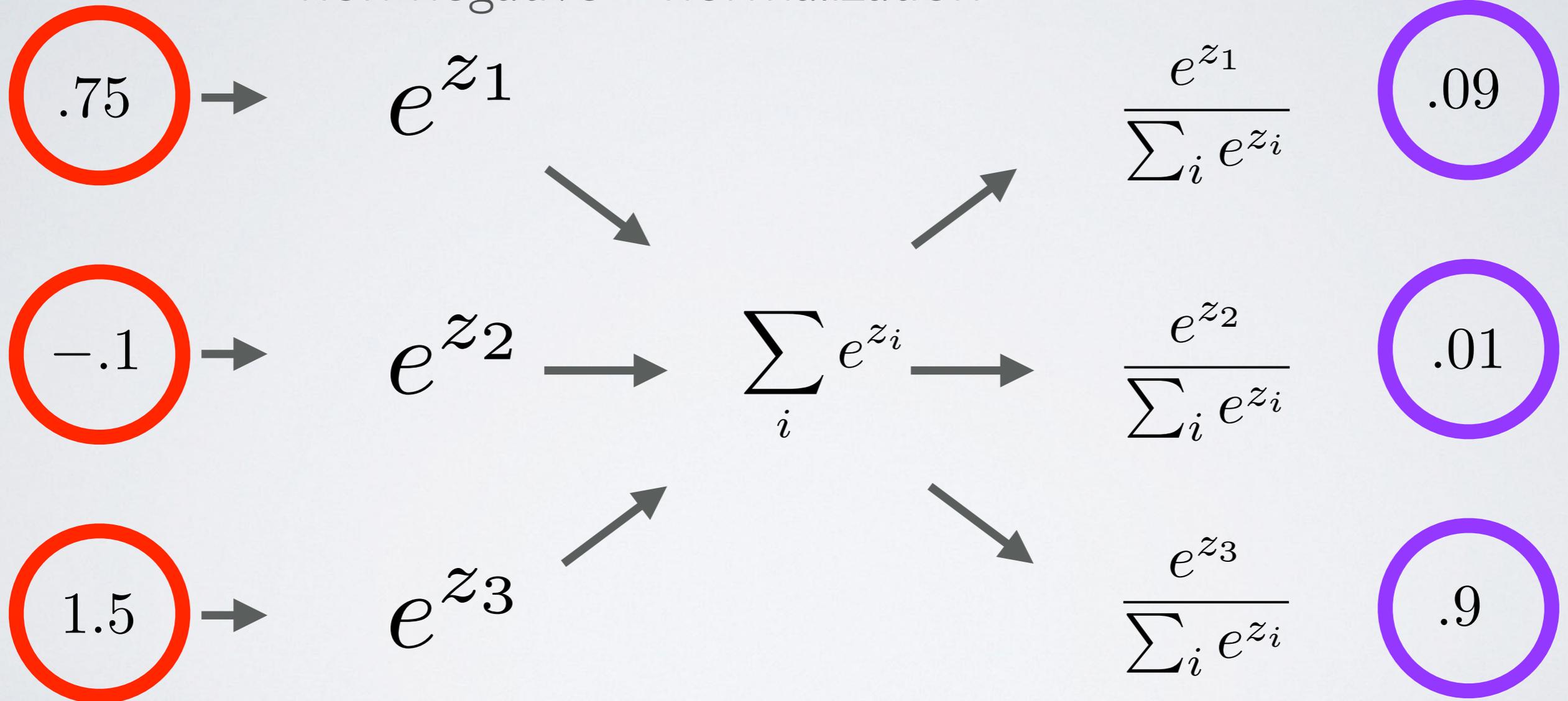


MULTI-CLASS NETWORK



SOFTMAX

make non-negative compute normalization probabilities



CROSS-ENTROPY LOSS

target label



y_i

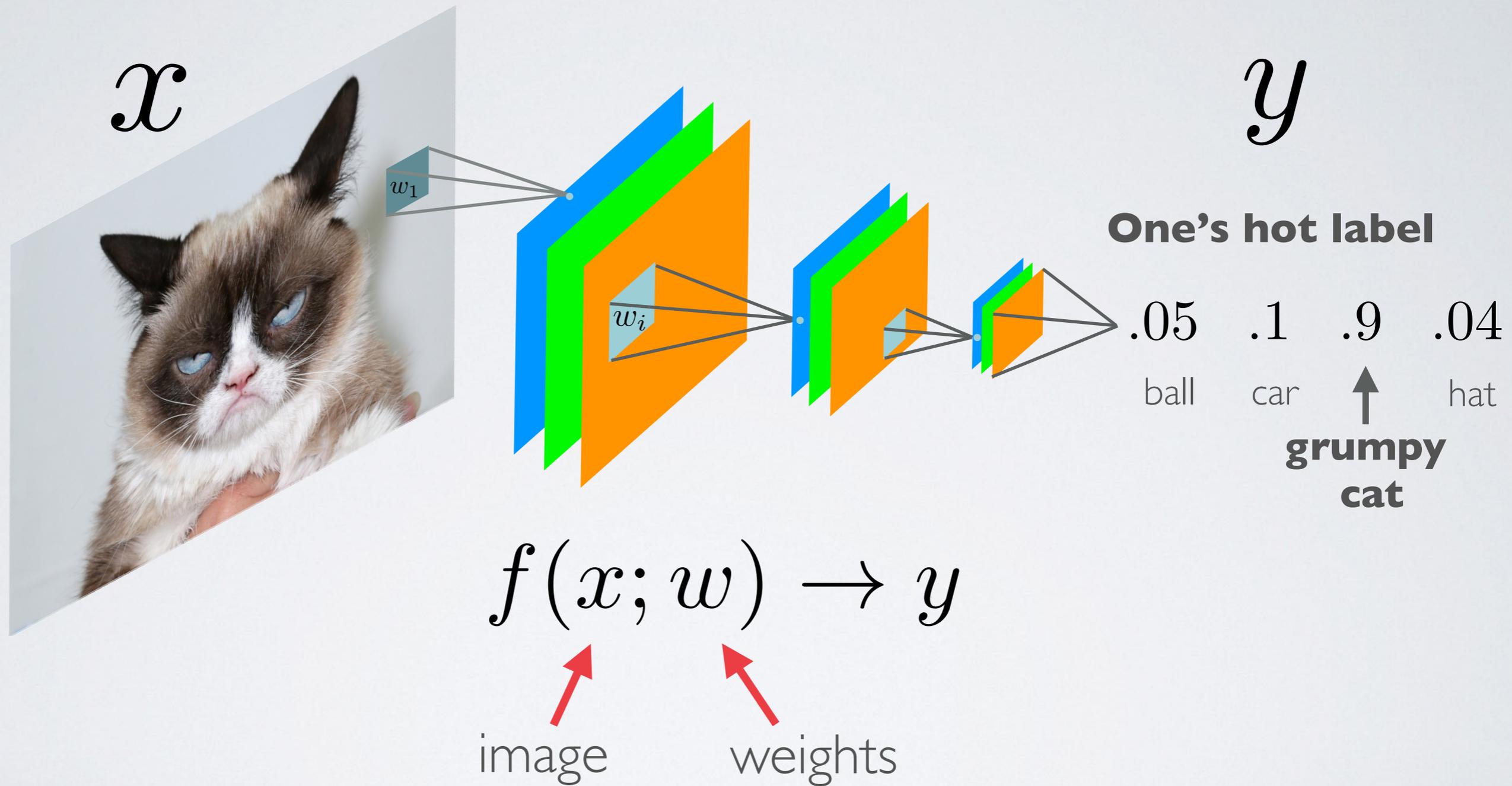
“cross entropy loss”
is negative log likelihood

$$\ell(x, y) = - \sum_k y^k \log(x^k)$$

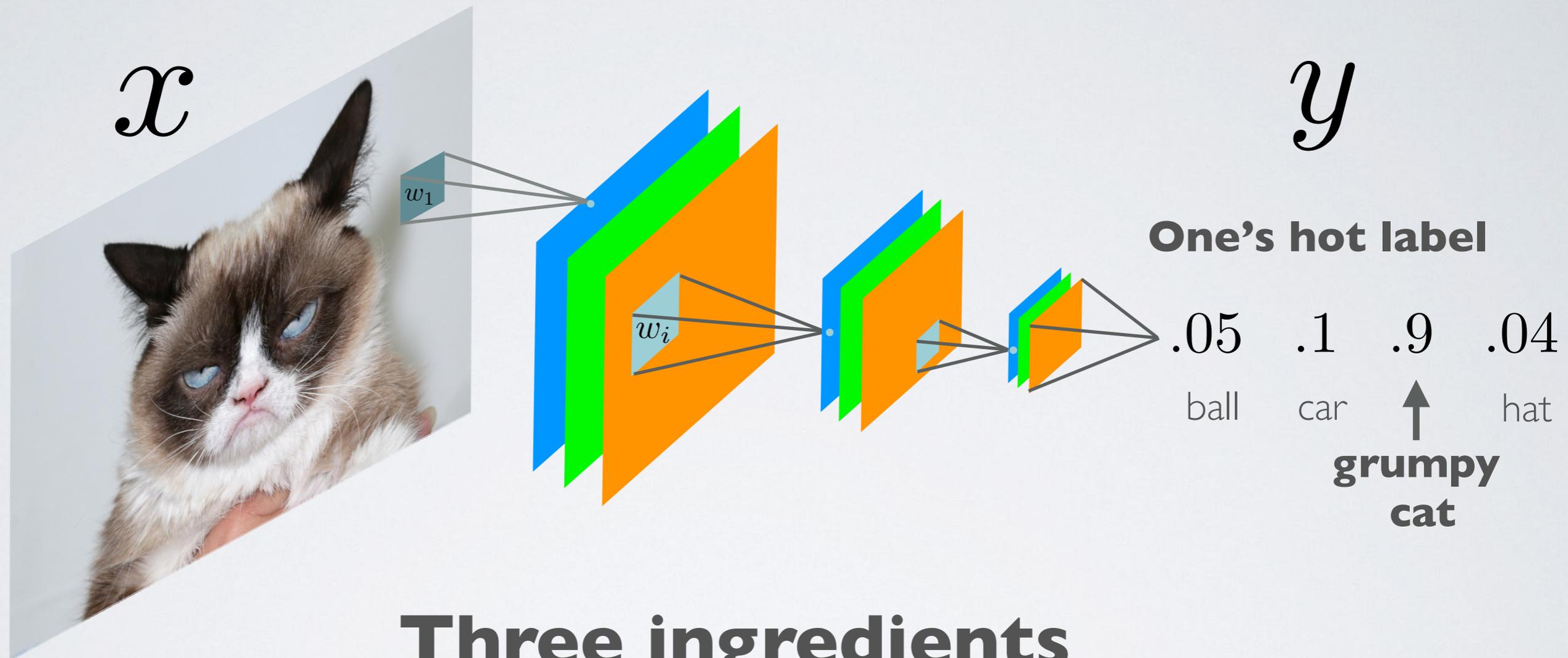
NLL = $-\log(.9)$

BUILDING MORE COMPLEX NETS

CONVOLUTIONAL NET



CONVOLUTIONAL NET



Three ingredients

- convolutions
- non-linearities
- skip connections

CONVOLUTION

**Input
image**

x

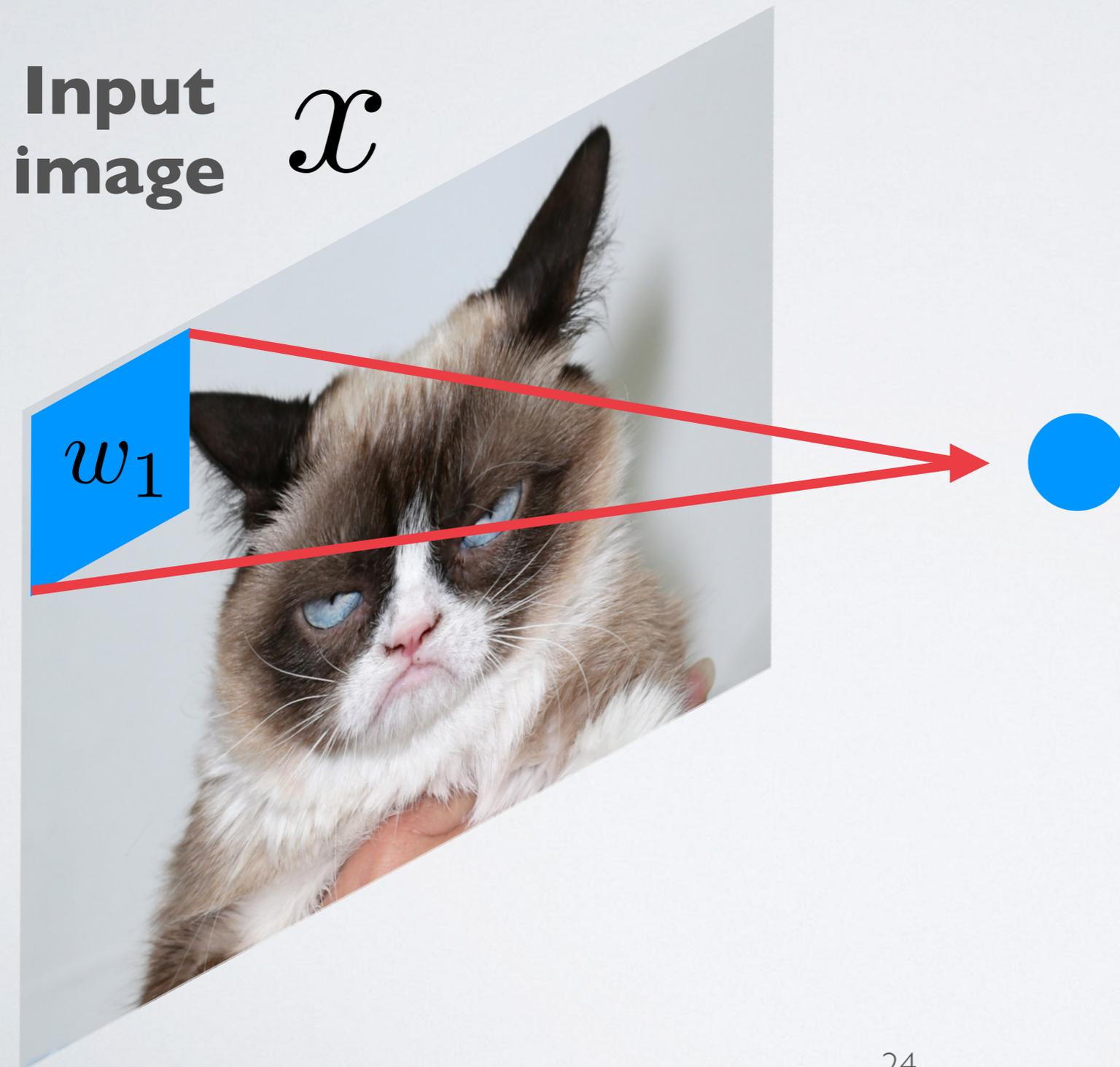


Convolutional filter

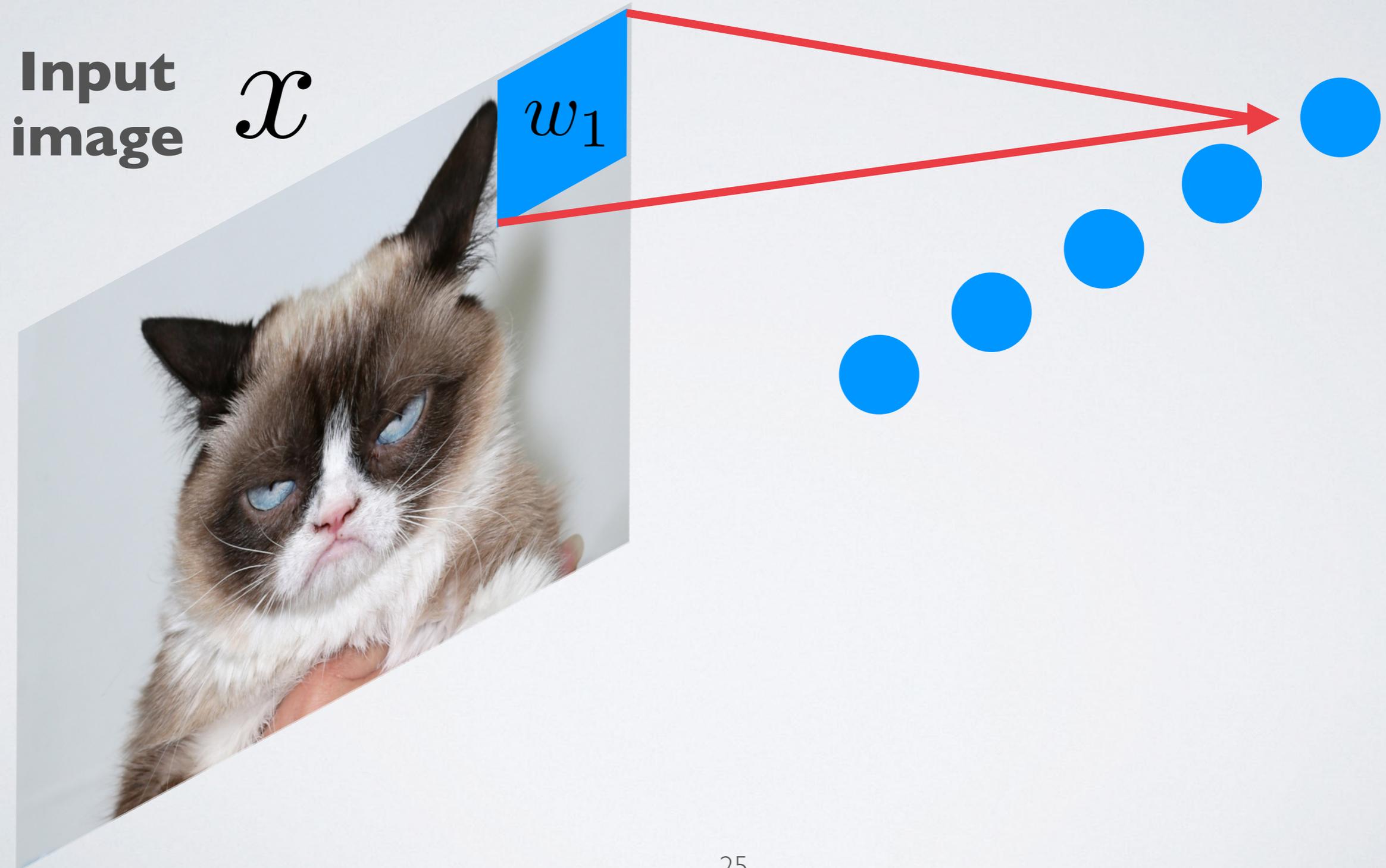
w_1

5	3	8
2	1	7
1	9	4

CONVOLUTION



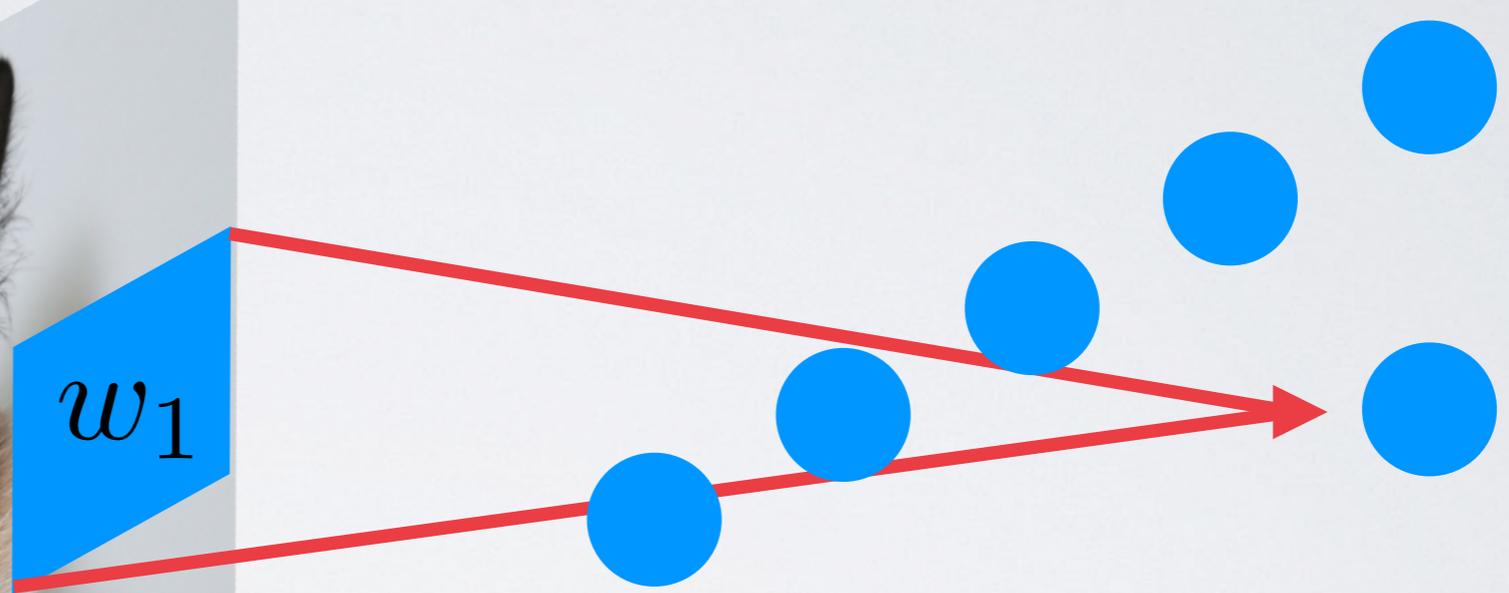
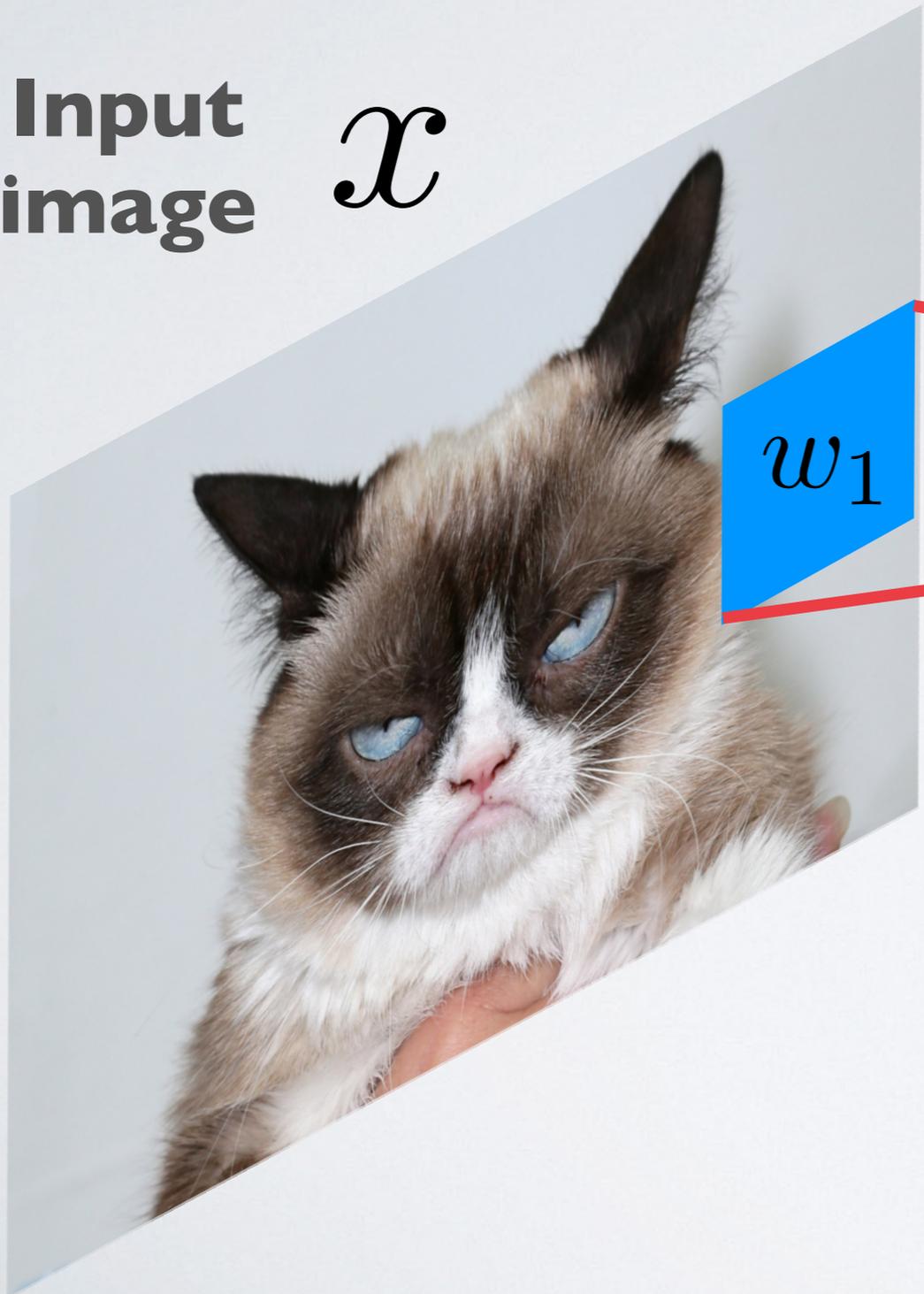
CONVOLUTION



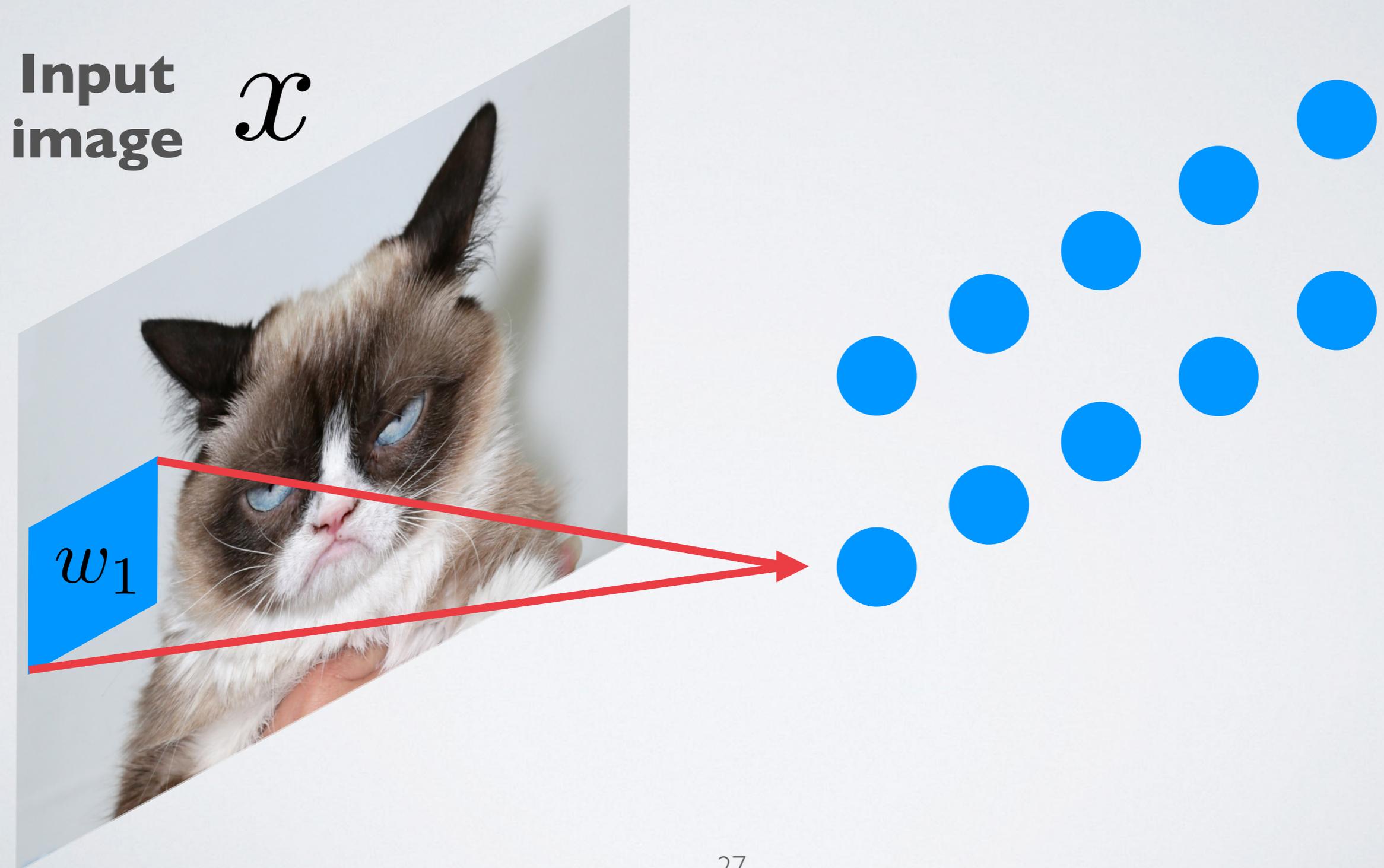
CONVOLUTION

Input
image

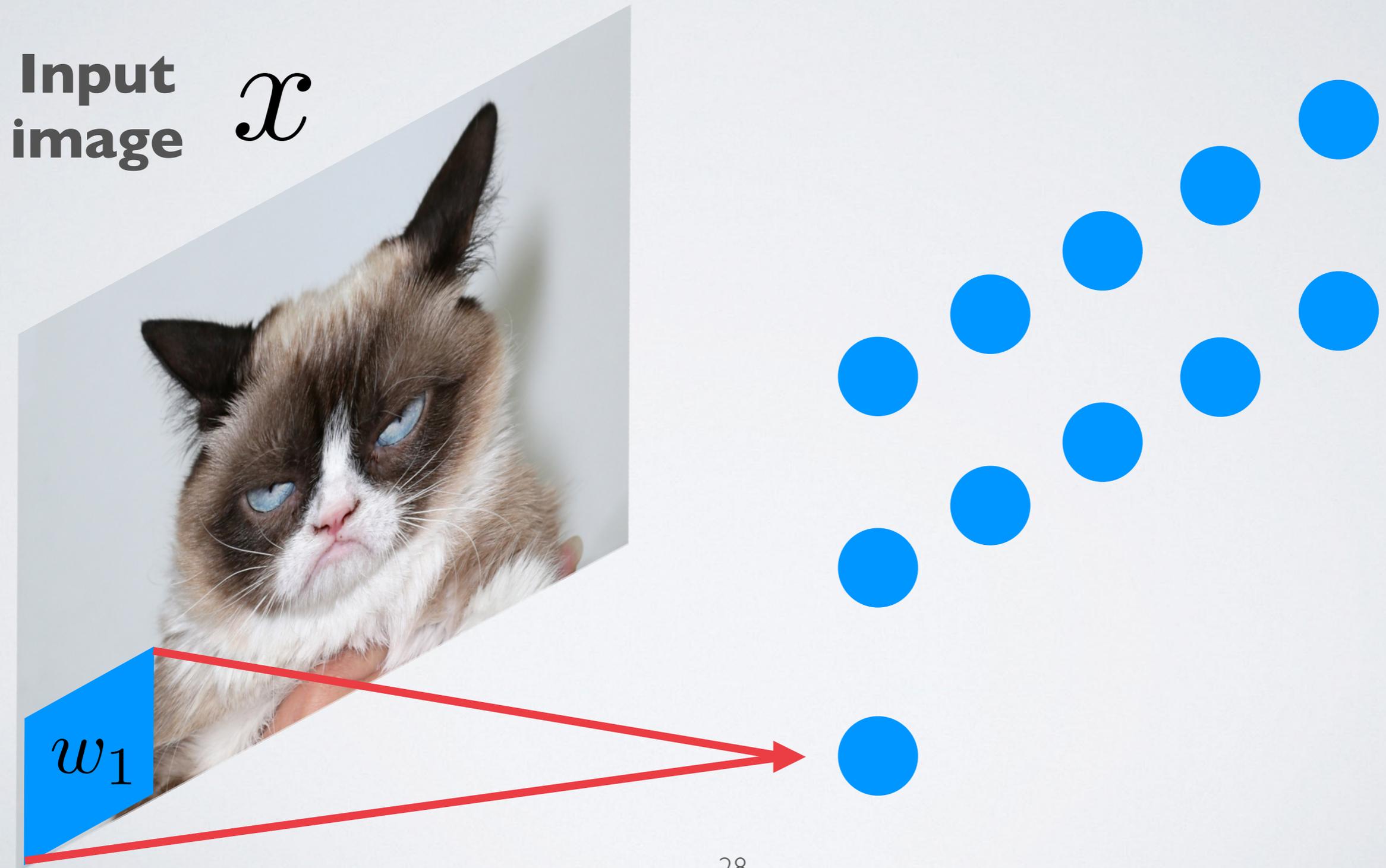
\mathcal{X}



CONVOLUTION



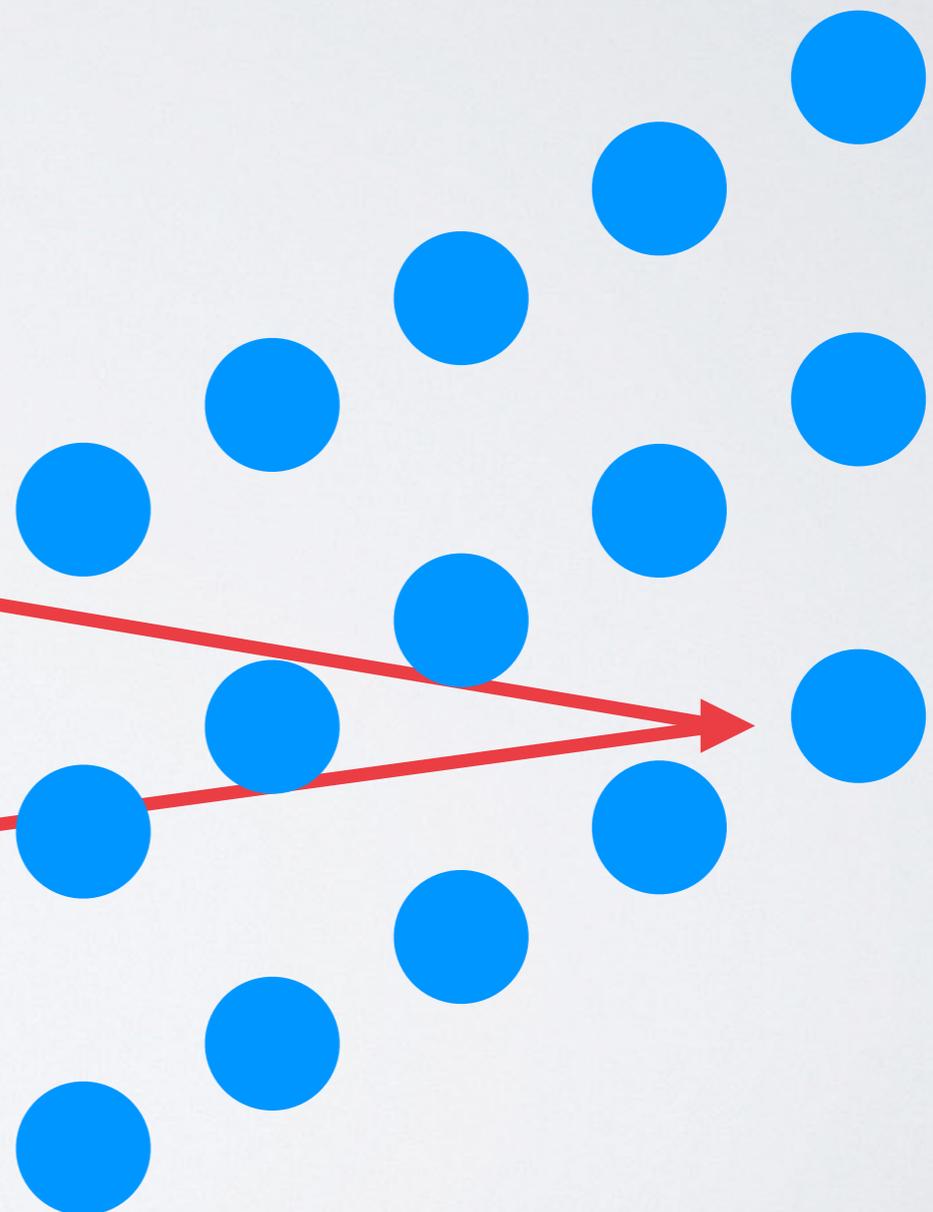
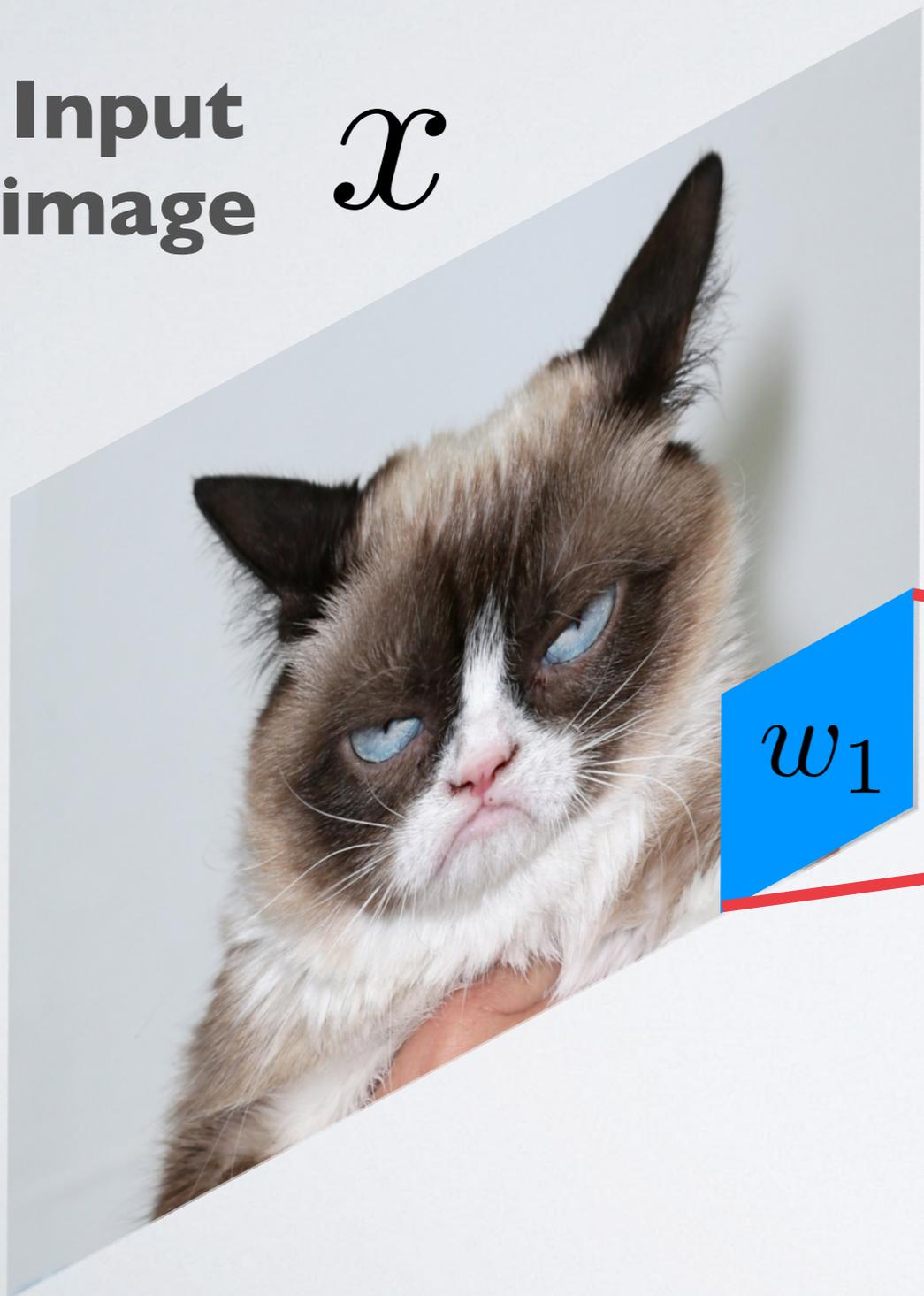
CONVOLUTION



CONVOLUTION

Input
image

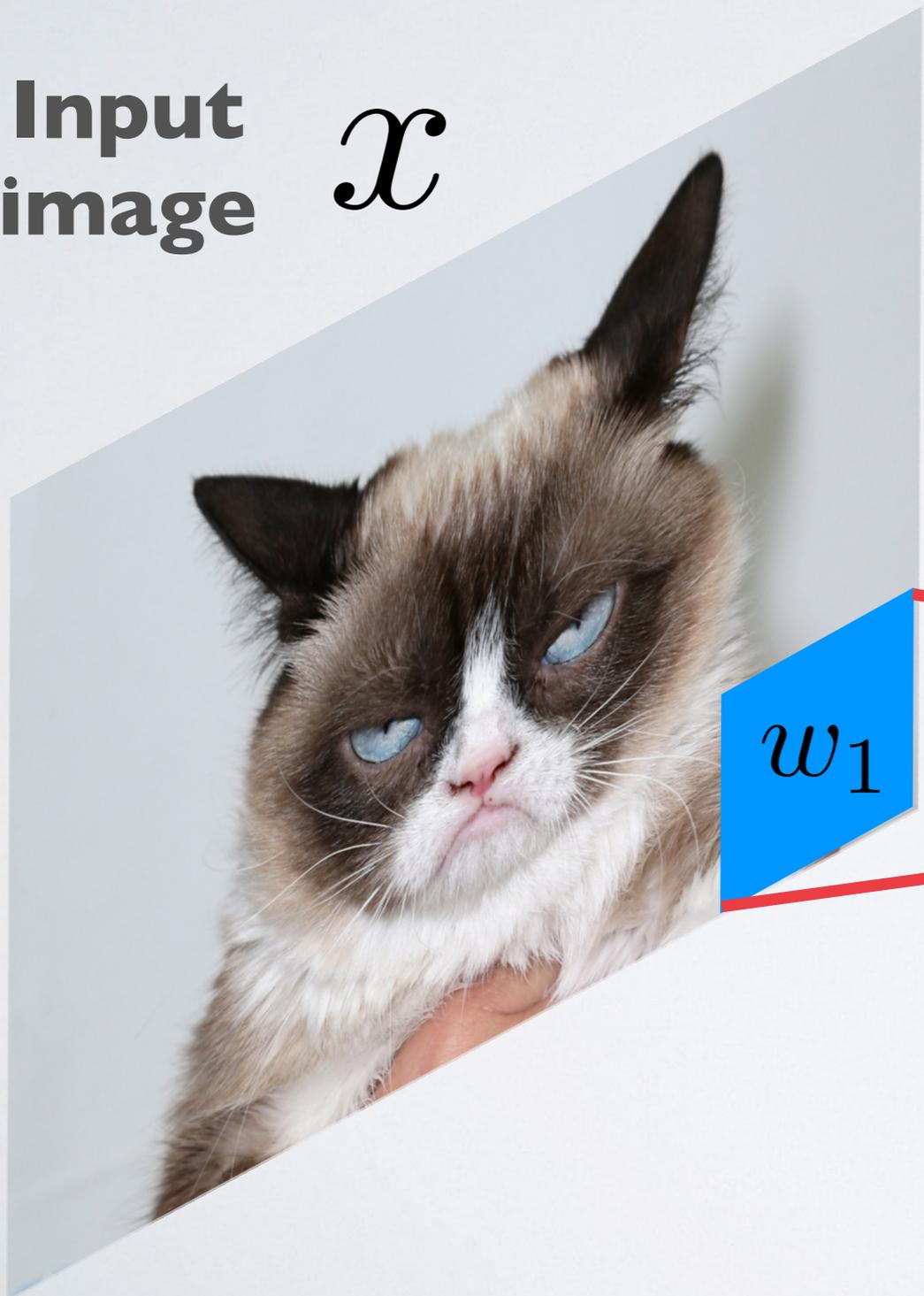
\mathcal{X}



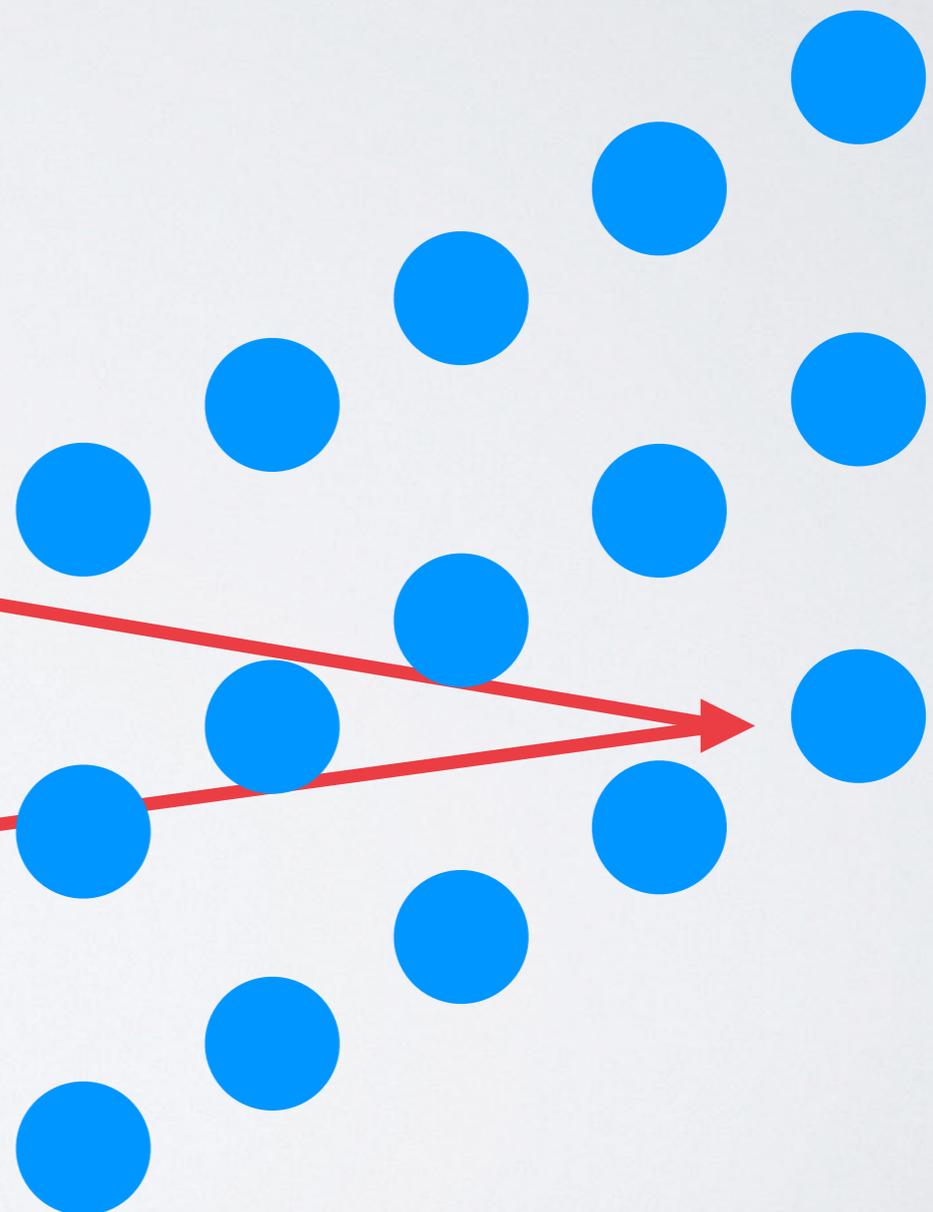
CONVOLUTION

Input
image

\mathcal{X}

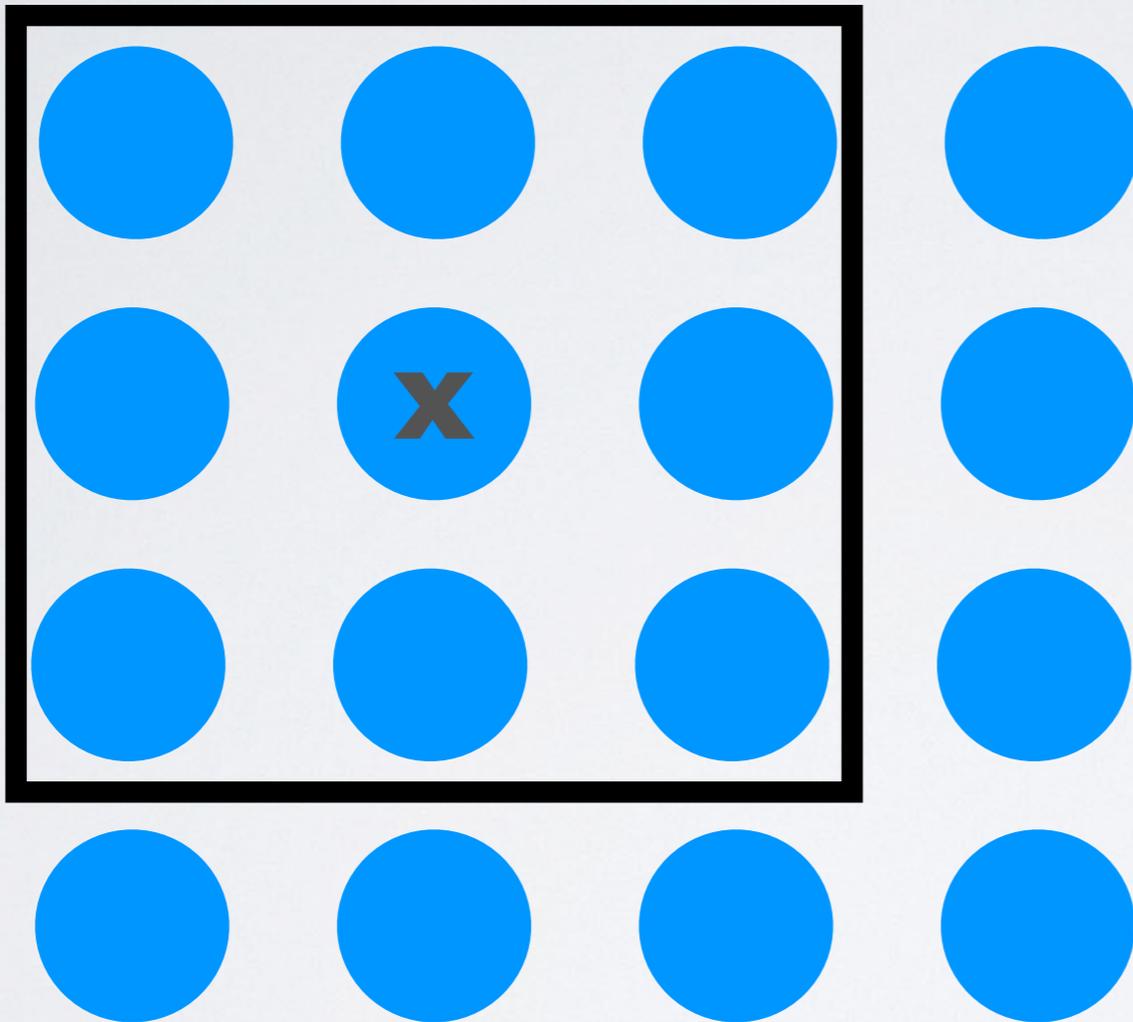


“Feature map”



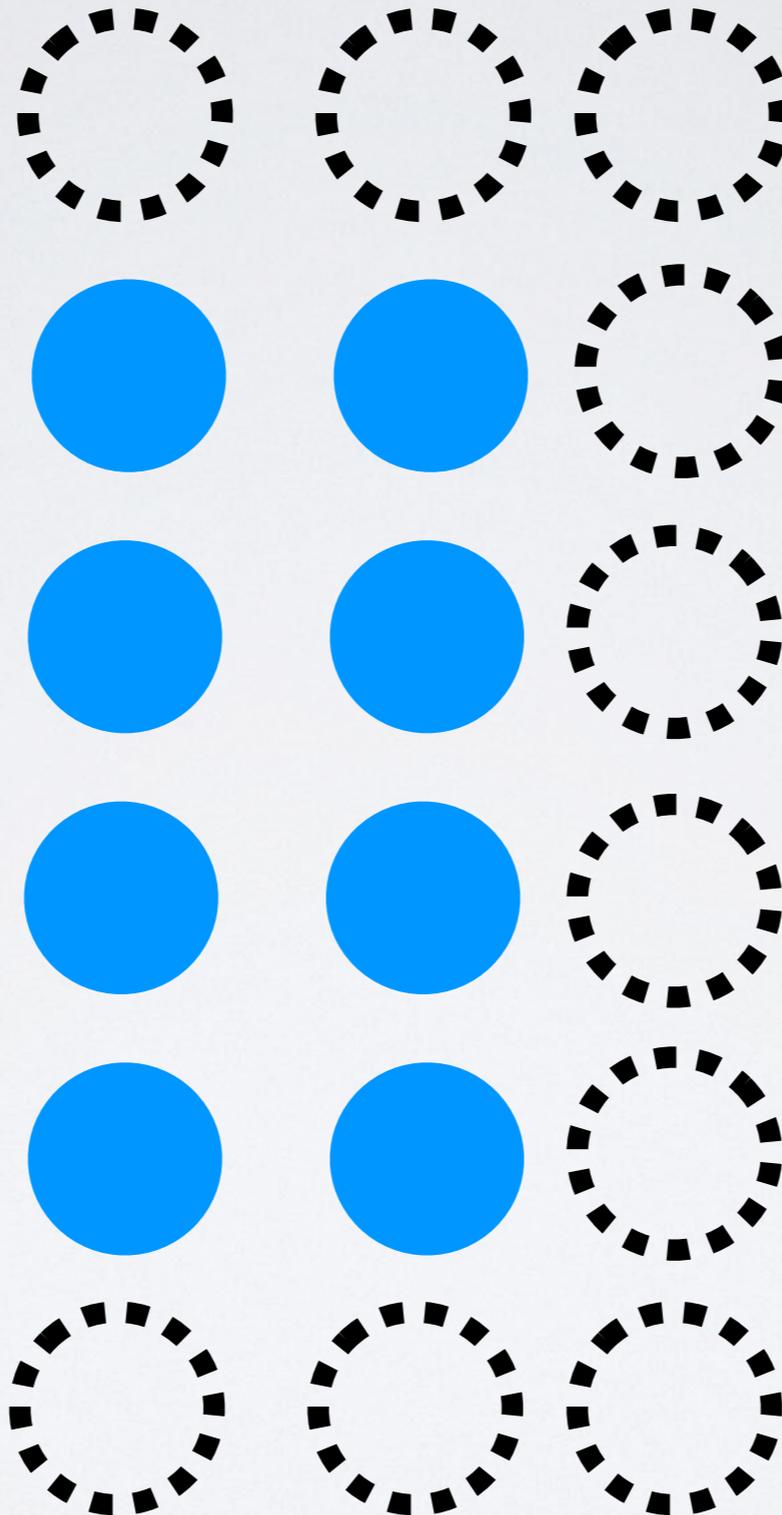
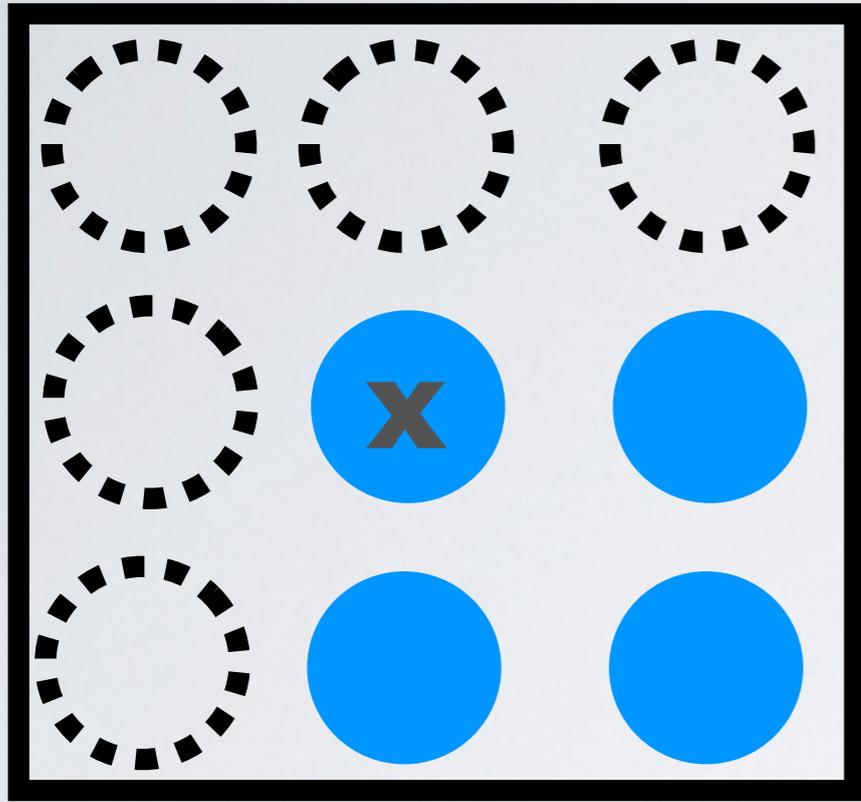
VARIANTS OF CONV

PADDING



**3x3 conv
produces
feature maps
that are 2 units
smaller**

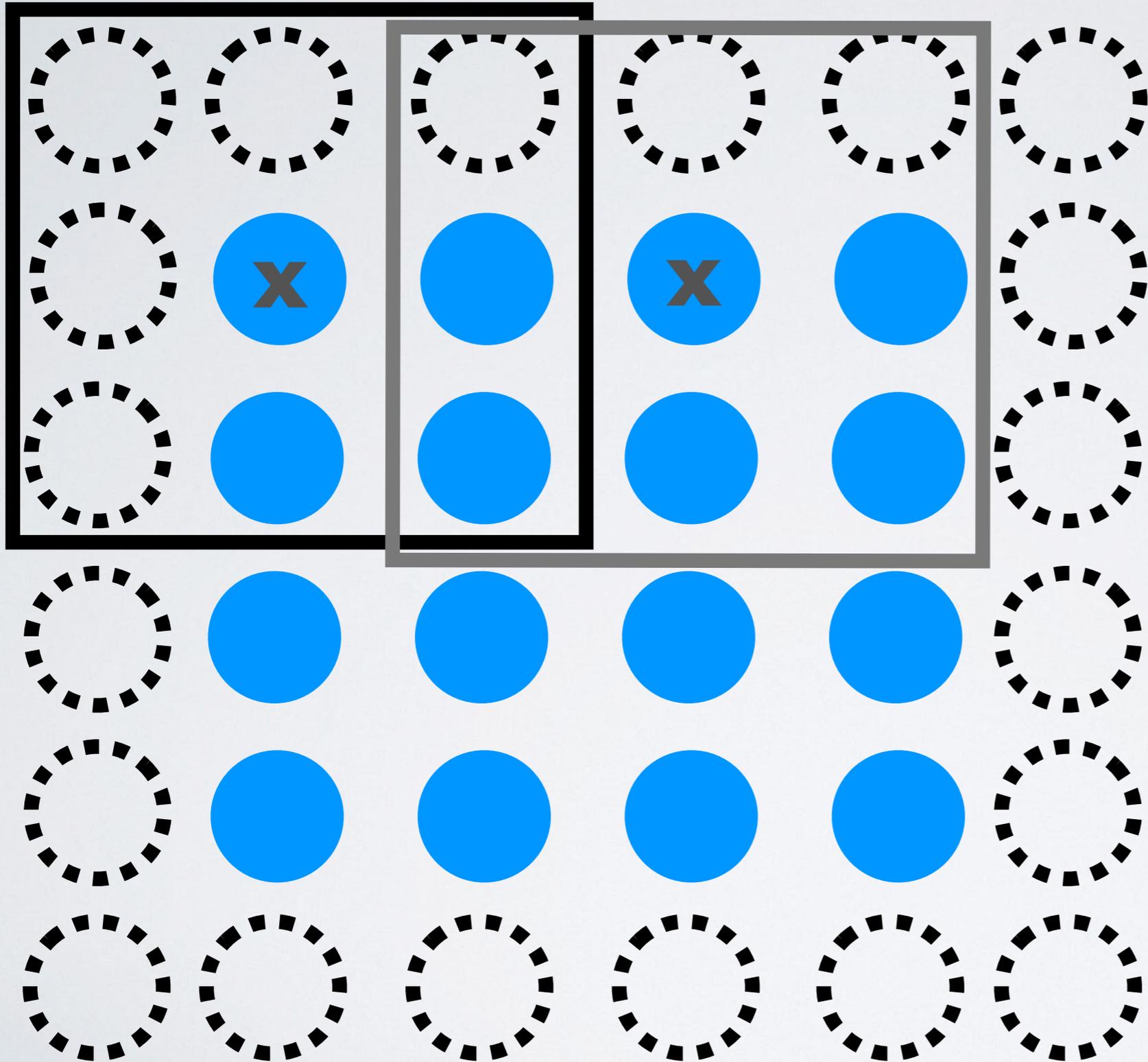
PADDING



padding=1

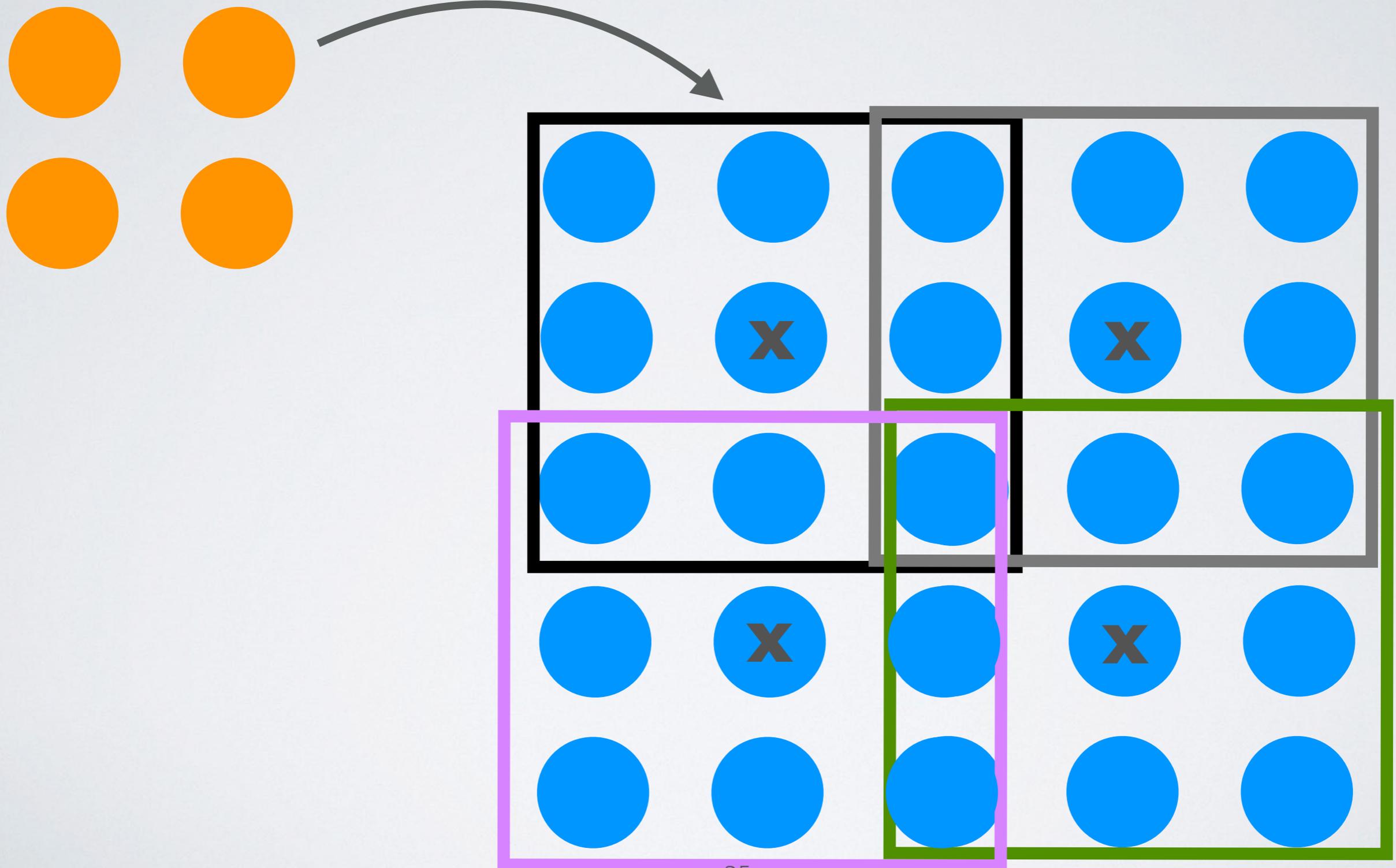
Output is same size
as input

STRIDE



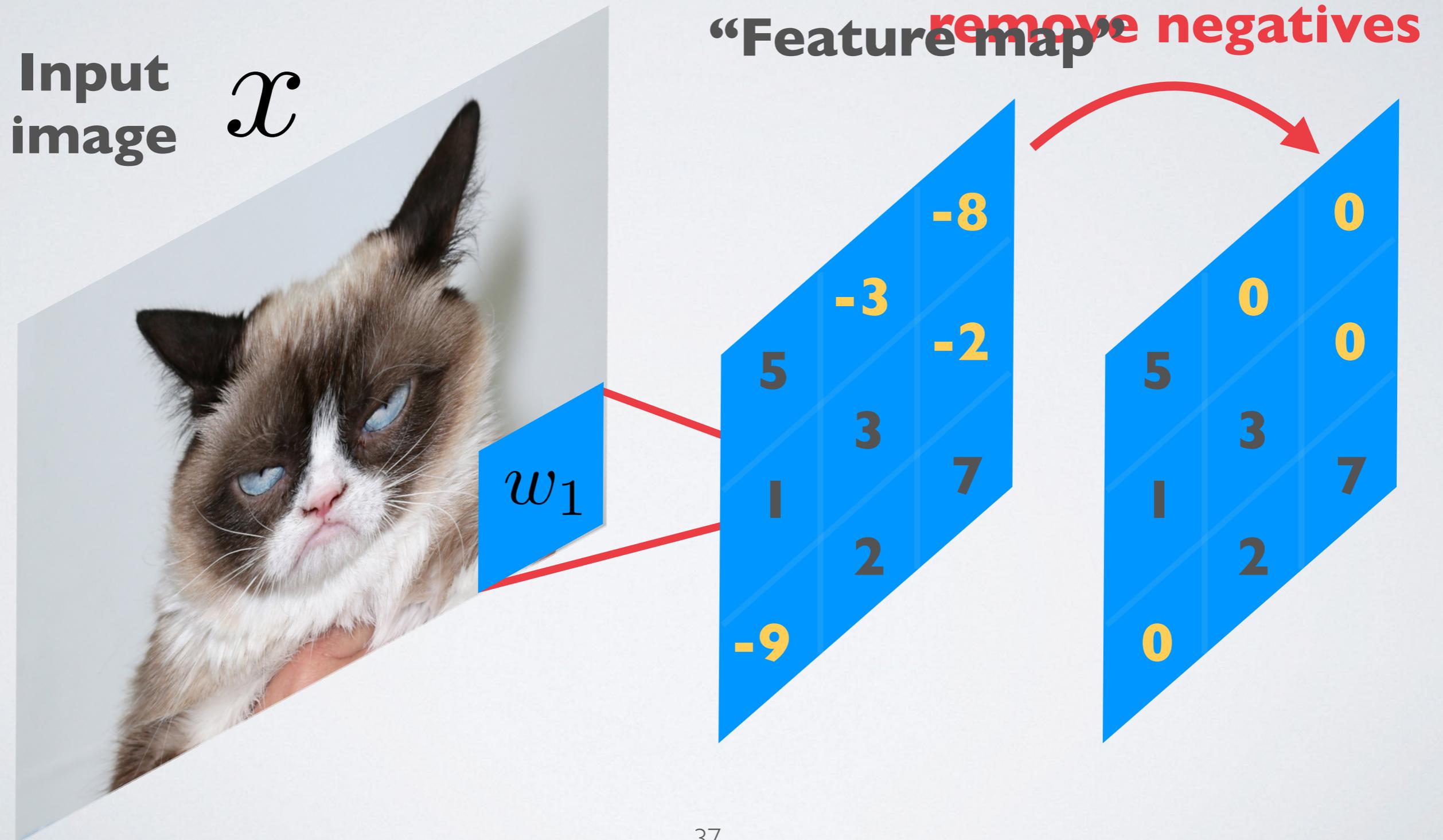
**stride=2 cuts
resolution by
half**

TRANSPOSE / DECONVOLUTION



RELU

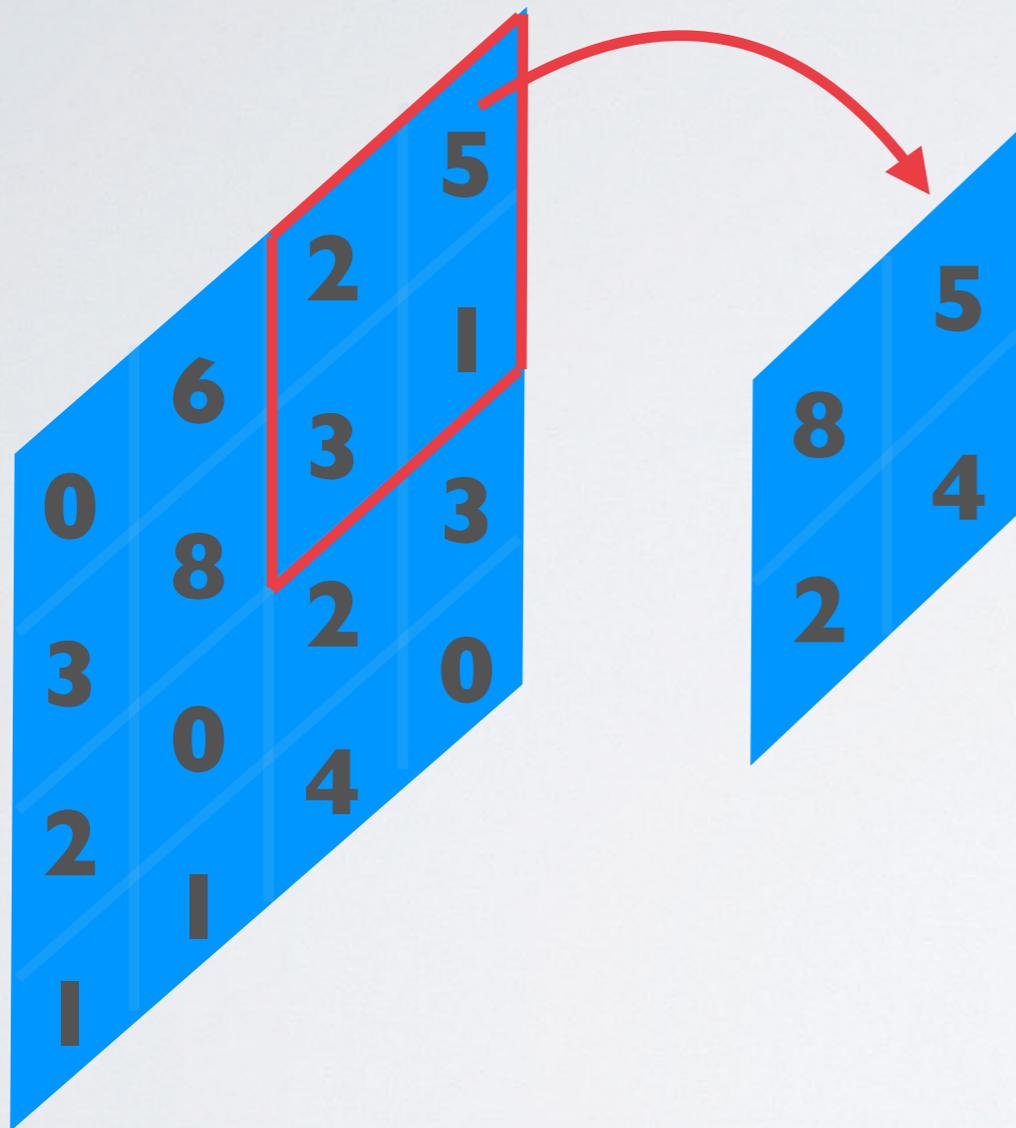
RELU NONLINEARITY



POOLING

POOLING

moving to lower resolutions



Max pooling

Average pooling

Strided convolution
(aka downsampling)

THAT'S ALL YOU NEED!

