

DS 503: Advanced Data Analytics

# Lecture 14: Into to Neural Nets

Based on MMDS Ch 14 and  
CS771 slides by Dr. Piyush Rai

Instructor: Dr. Gagan Gupta

# Example (Learning Non-Linear Boundaries)

Example 1: We have length 4 bit-vectors of the form  $[x_1 \ x_2 \ x_3 \ x_4]$ .  
The output  $y$  is 1 when there are 2 consecutive 1s.

Training Ex

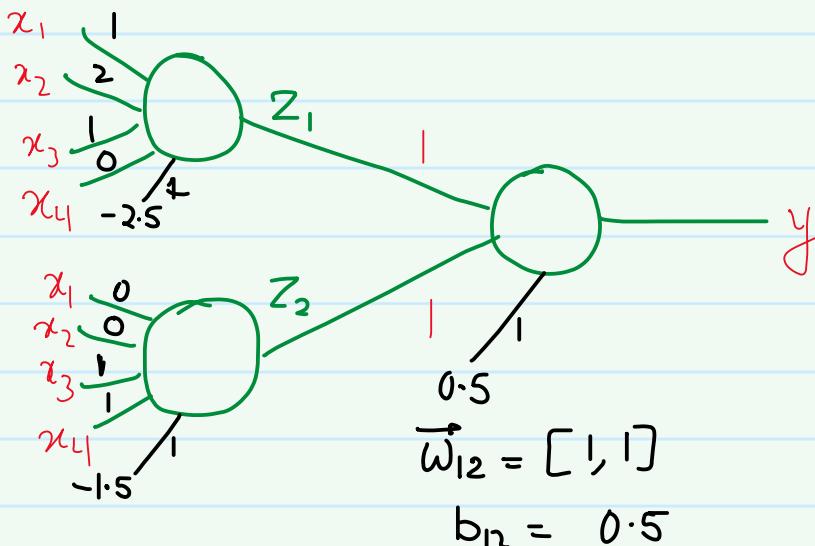
$$\begin{bmatrix} 0 & 1 & 1 & 0 \end{bmatrix}, 1$$
$$\begin{bmatrix} 0 & 0 & 1 & 0 \end{bmatrix}, 0$$

$$\begin{bmatrix} 1 & 1 & 0 & 0 \end{bmatrix}, 1$$
$$\begin{bmatrix} 1 & 0 & 0 & 1 \end{bmatrix}, 0$$

and so on

$$\vec{\omega}_1 = [1, 2, 1, 0]$$
$$b_1 = -2.5$$

$$\vec{\omega}_{21} = [0, 0, 1, 1]$$
$$b_{21} = -1.5$$



Each node is a perceptron.

Q. When can  $z_1$  be 1?

$$x_1 + 2x_2 + x_3 - 2.5 > 0$$

$$\Rightarrow [x_2 > 0] \times [x_1 \text{ or } x_3 > 0]$$

Q. When can  $z_2$  be 1?

Q. When can  $y$  be 1?

$$[z_1 > 0] \text{ or } [z_2 > 0]$$

# Illustration: Neural Net with One Hidden Layer

- Each input  $\mathbf{x}_n$  transformed into several pre-activations using linear models

$$a_{nk} = \mathbf{w}_k^\top \mathbf{x}_n = \sum_{d=1}^D w_{dk} x_{nd}$$

- Nonlinear activation applied on each pre-act.

$$h_{nk} = g(a_{nk})$$

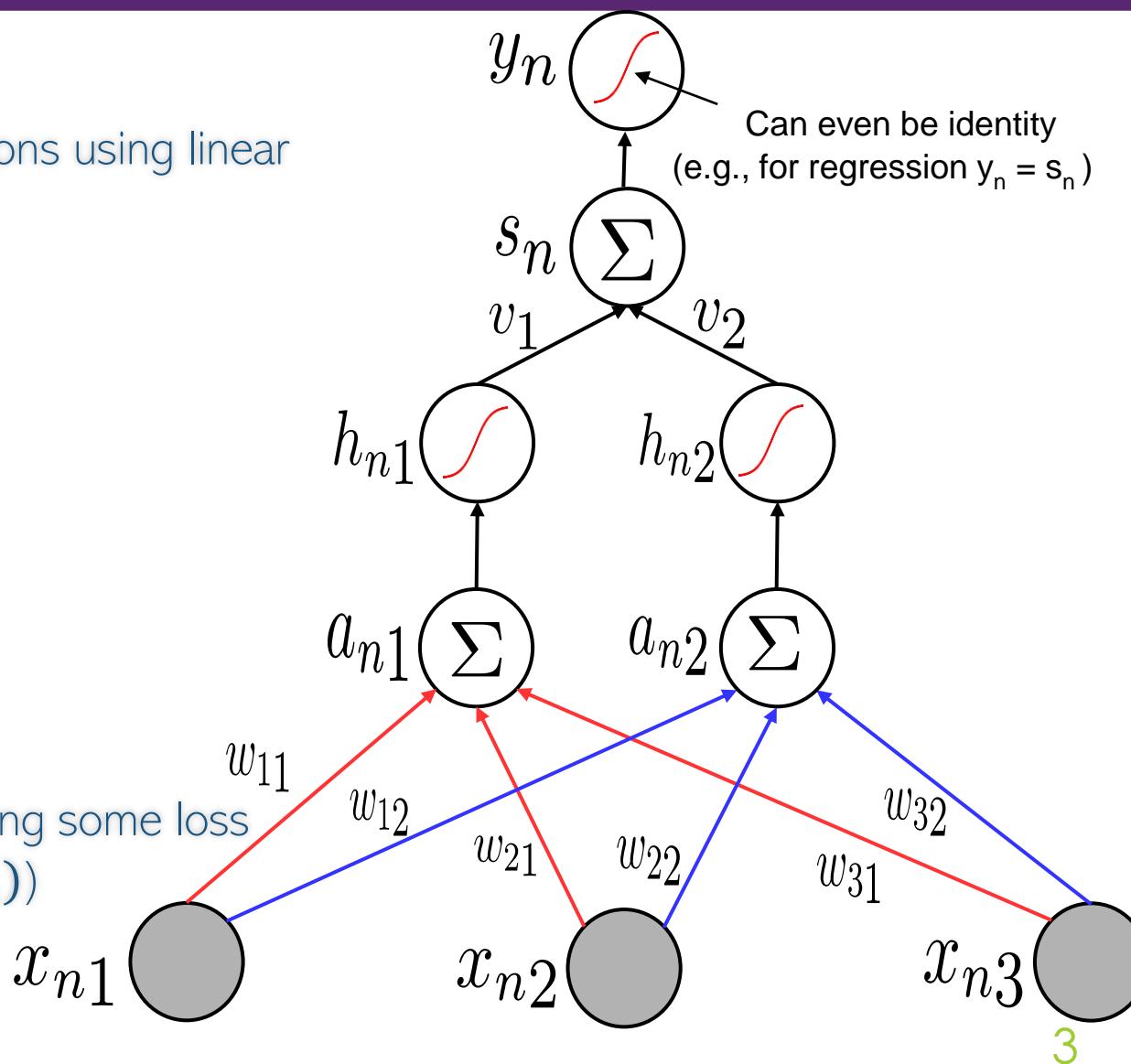
- Linear model learned on the new “features”  $\mathbf{h}_n$

$$s_n = \mathbf{v}^\top \mathbf{h}_n = \sum_{k=1}^K v_k h_{nk}$$

- Finally, output is produced as  $y = o(s_n)$

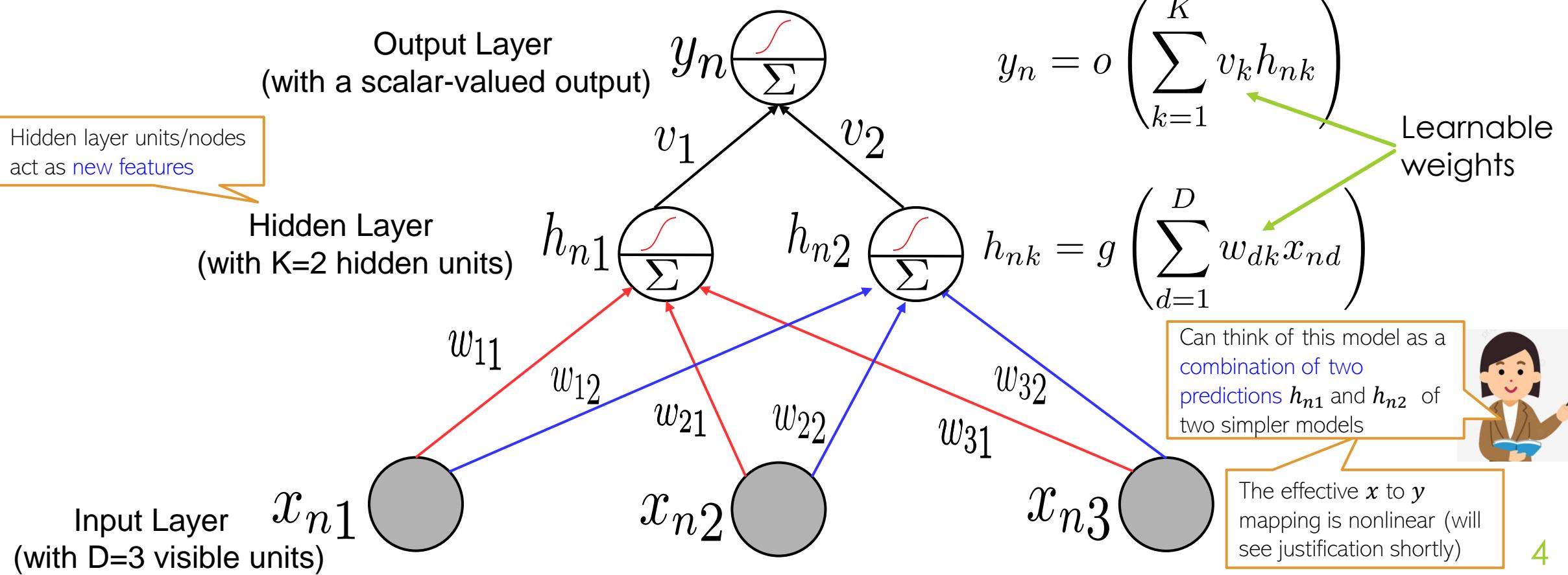
- Unknowns  $(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K, \mathbf{v})$  learned by minimizing some loss function, for example  $\mathcal{L}(\mathbf{W}, \mathbf{v}) = \sum_{n=1}^N \ell(y_n, o(s_n))$

(squared, logistic, softmax, etc)



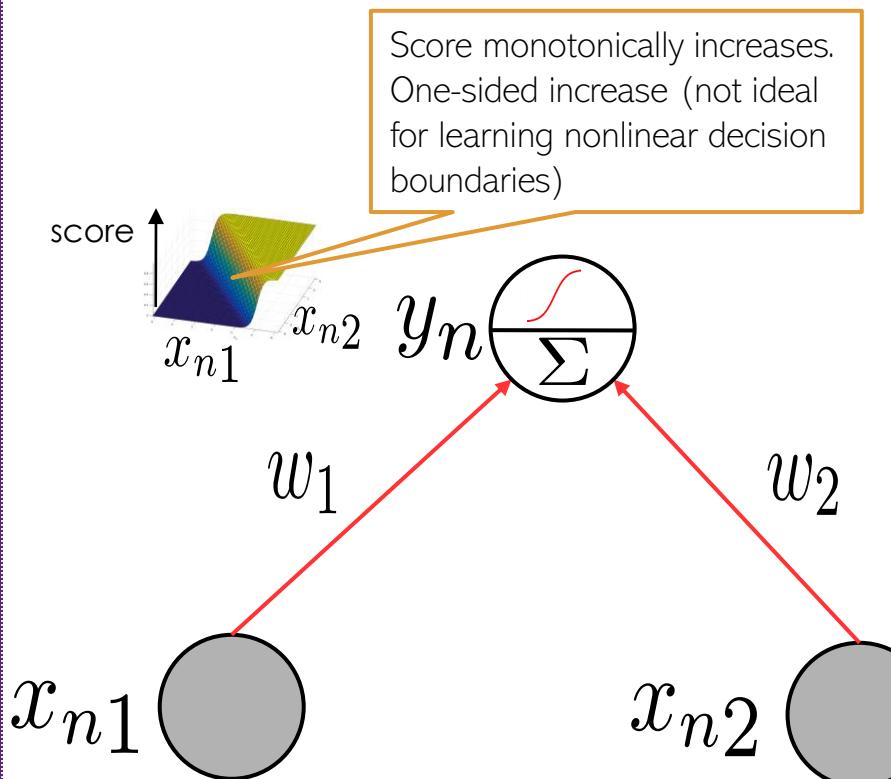
# Neural Networks: Multi-layer Perceptron (MLP)

- An MLP consists of an **input layer**, an **output layer**, and **one or more hidden layers**

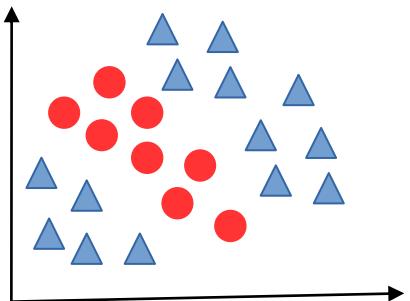


# MLP Can Learn Nonlin. Fn: A Brief Justification

An MLP can be seen as a composition of multiple linear models combined nonlinearly



Obtained by composing the two one-sided increasing score functions (using  $v_1 = 1$ , and  $v_2 = -1$  to "flip" the second one before adding). This can now learn nonlinear decision boundary

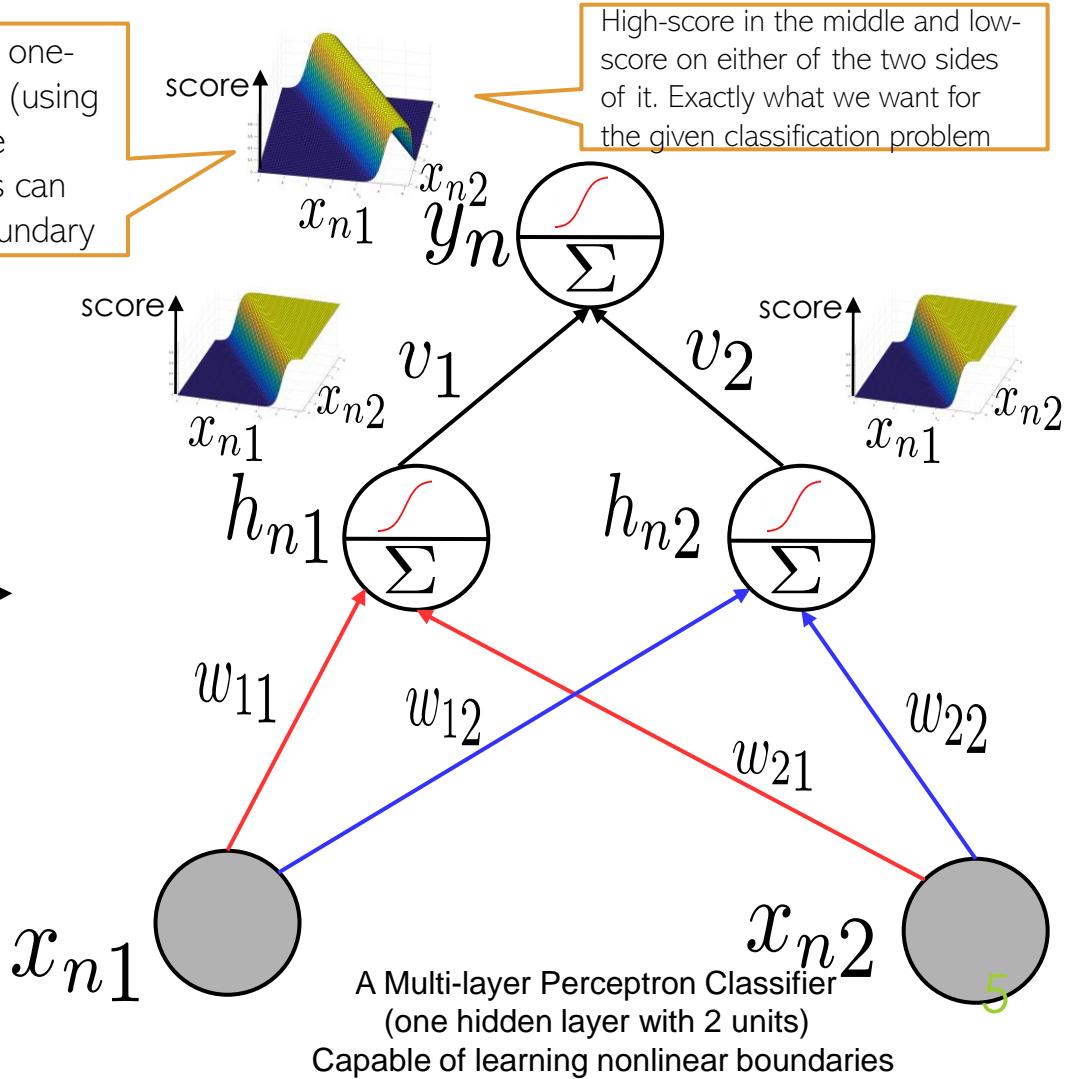


A nonlinear classification problem

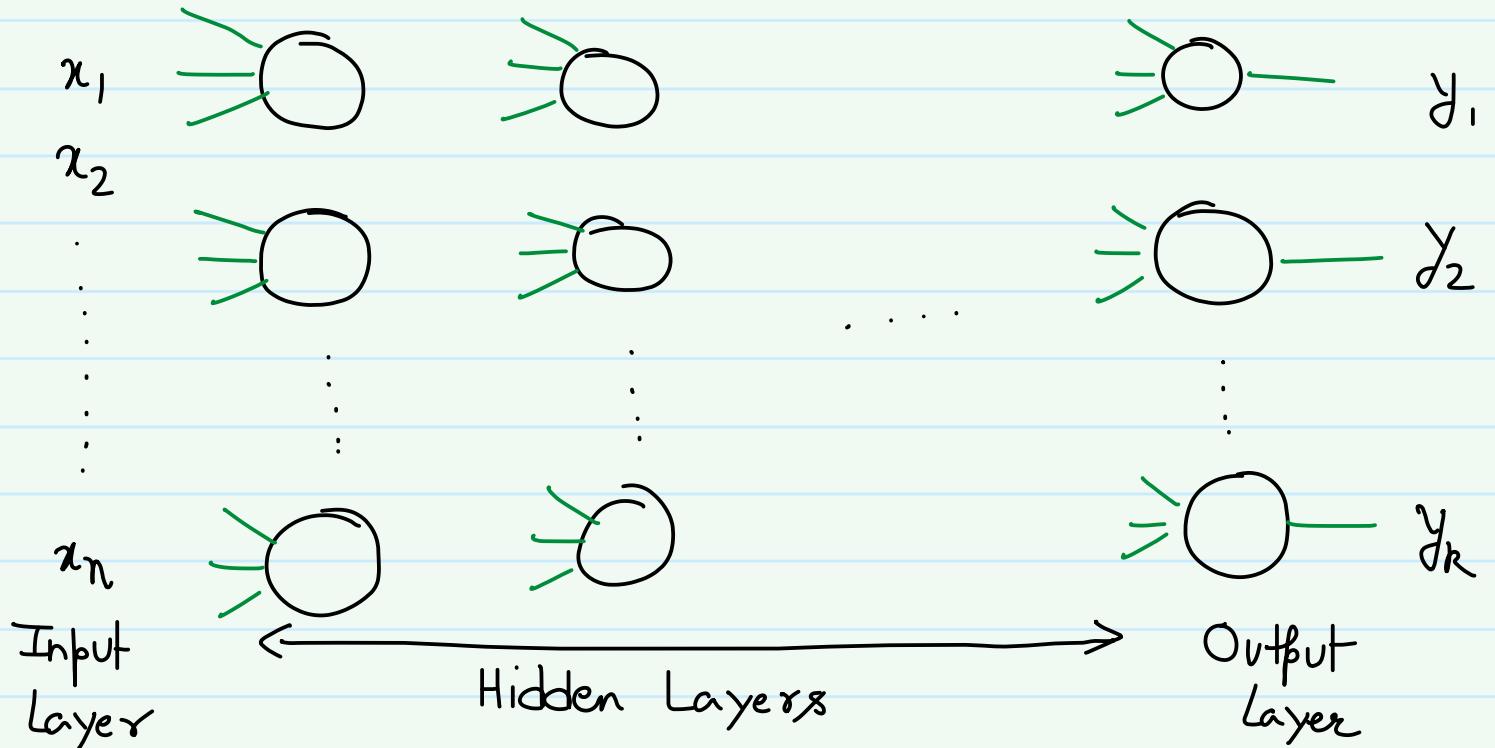


Standard Single "Perceptron" Classifier (no hidden units)

A single hidden layer MLP with sufficiently large number of hidden units can approximate any function (Hornik, 1991)



# Neural Net Designs



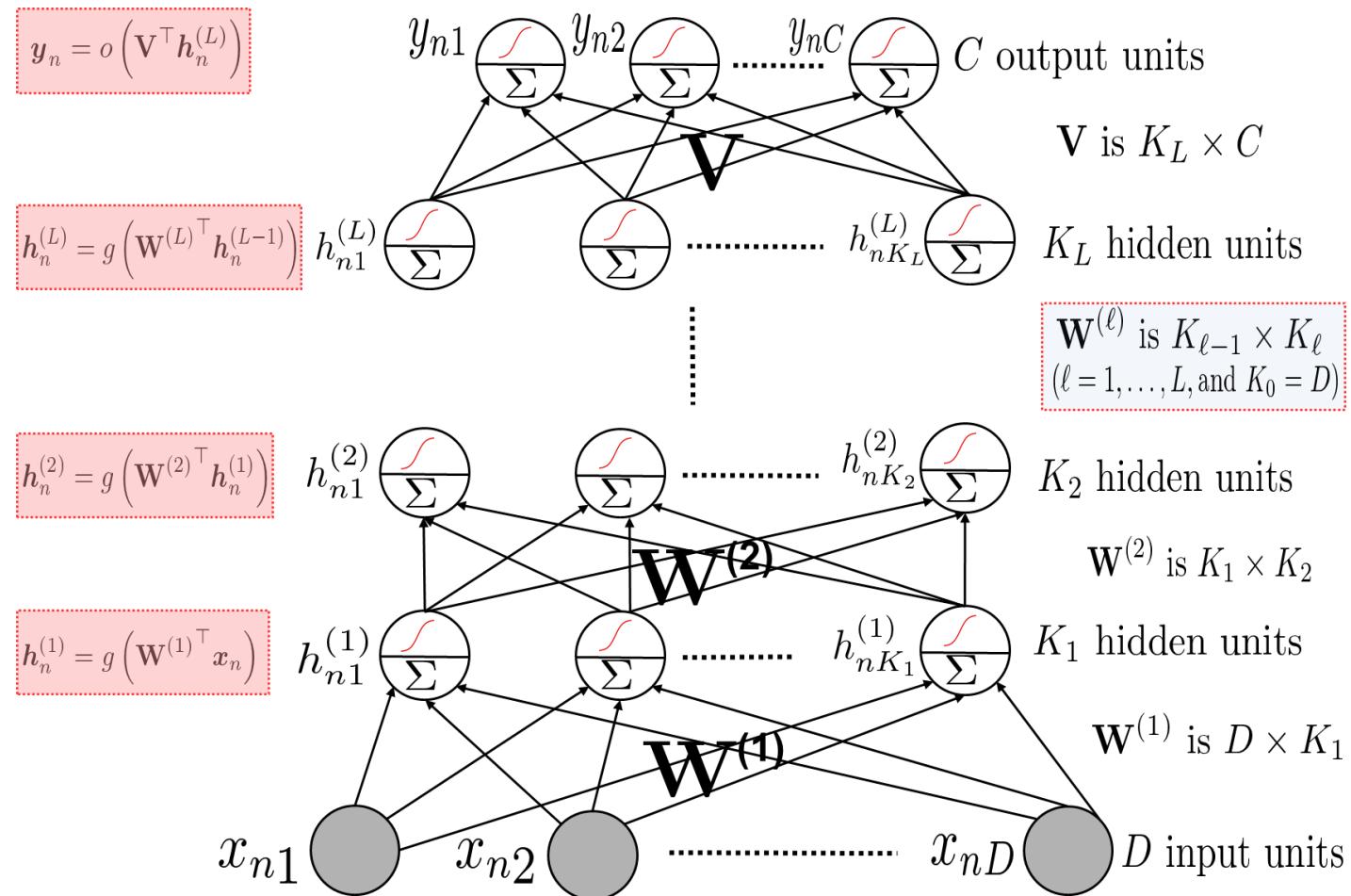
Output of each node depends on

$$\sum x_i w_i$$

where the sum is over all inputs  $x_i$  &  $w_i$  is the weight of the input.

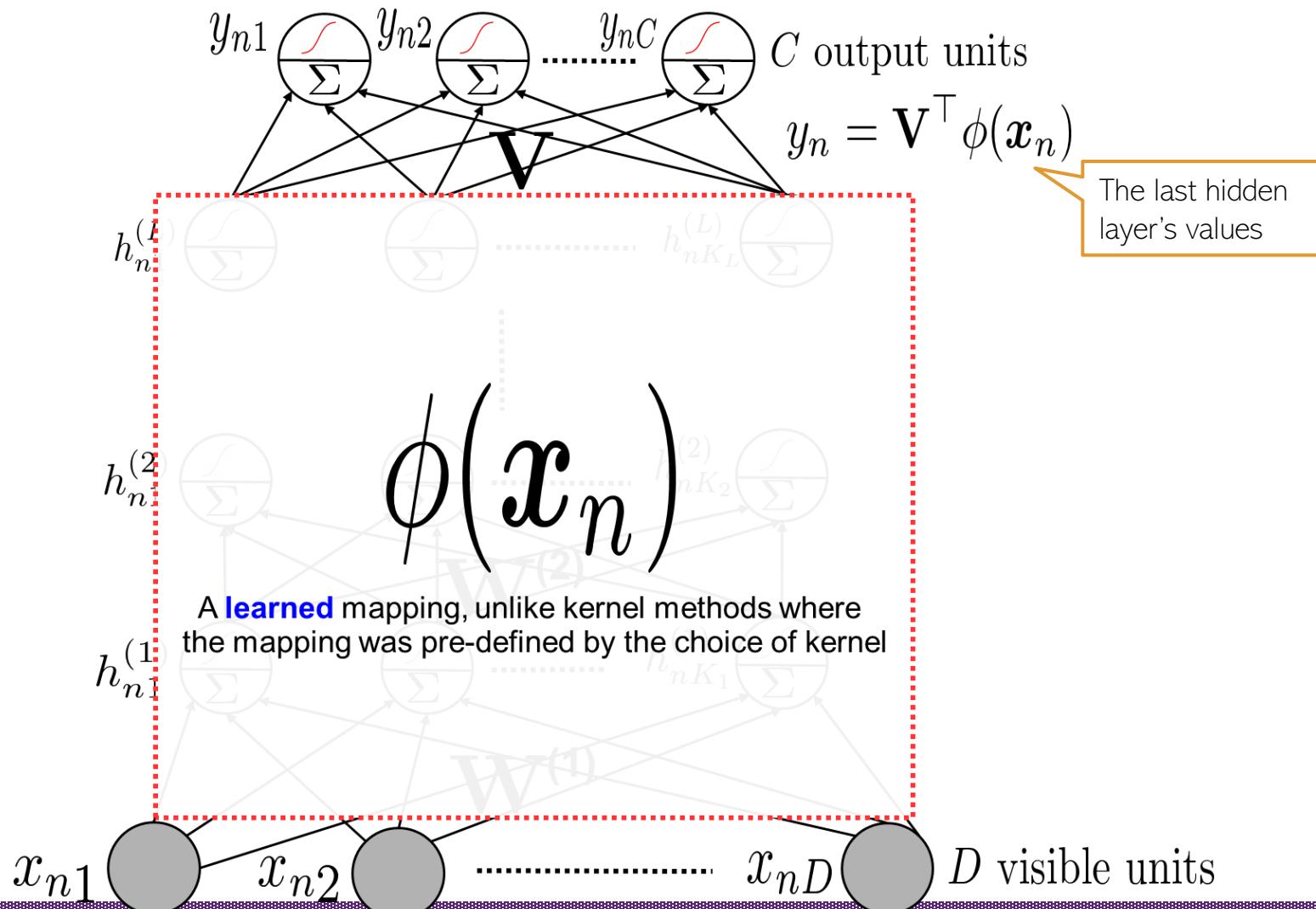
# Multiple Hidden Layers (One/Multiple Outputs)

- Most general case: Multiple hidden layers with (with same or different number of hidden nodes in each) and a scalar or vector-valued output



# Neural Nets are Feature Learners/Encoders

- Hidden layers can be seen as learning a feature rep.  $\phi(\mathbf{x}_n)$  for each input  $\mathbf{x}_n$



# Kernel Methods vs Neural Nets

- Prediction rule for a kernel method (e.g., kernel SVM)  
$$y = \sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \mathbf{x})$$
- This is analogous to a single hidden layer NN with fixed/pre-defined hidden nodes  $\{k(\mathbf{x}_n, \mathbf{x})\}_{n=1}^N$  and output weights  $\{\alpha_n\}_{n=1}^N$
- The prediction rule for a deep neural network

$$y = \sum_{k=1}^K v_k h_k$$

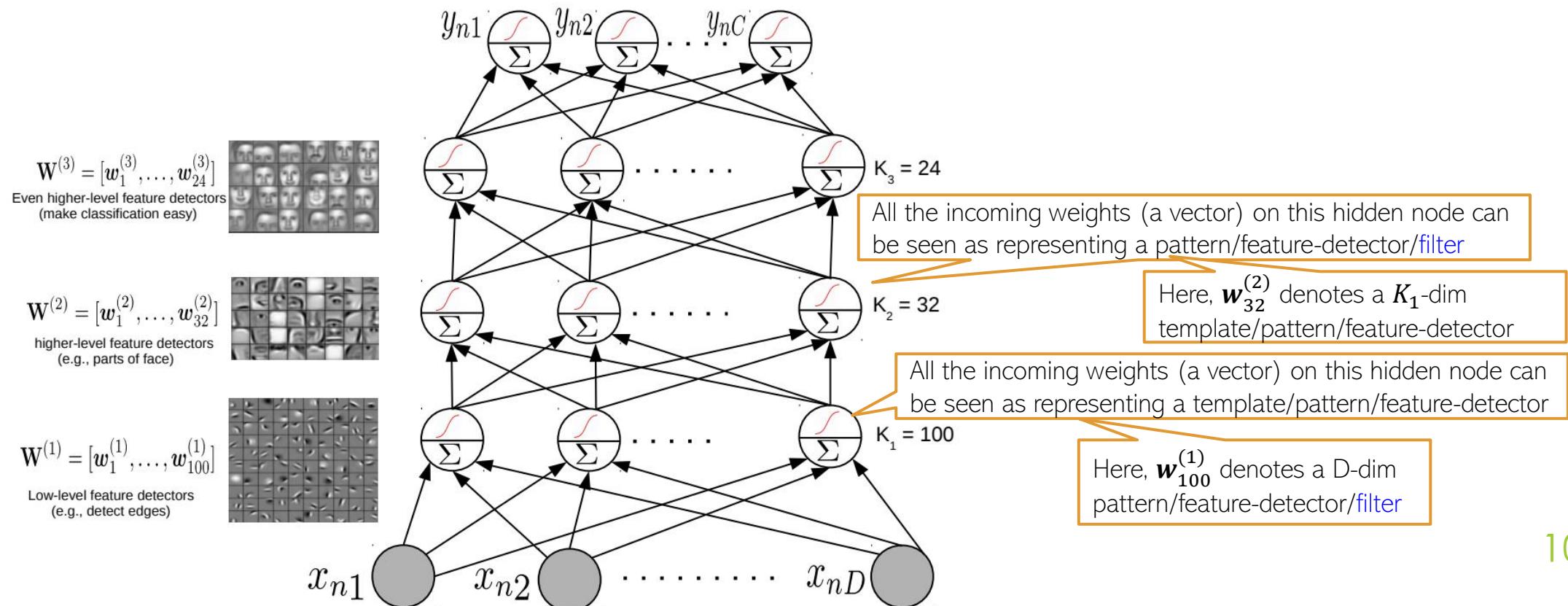
- Here, the  $h_k$ 's are learned from data (possibly after multiple layers of nonlinear transformations)
- Both kernel methods and deep NNs be seen as using nonlinear basis functions for making predictions. Kernel methods use fixed basis functions (defined by the kernel) whereas NN learns the basis functions adaptively from data

Also note that neural nets are faster than kernel methods at test time since kernel methods need to store the training examples at test time whereas neural nets do not



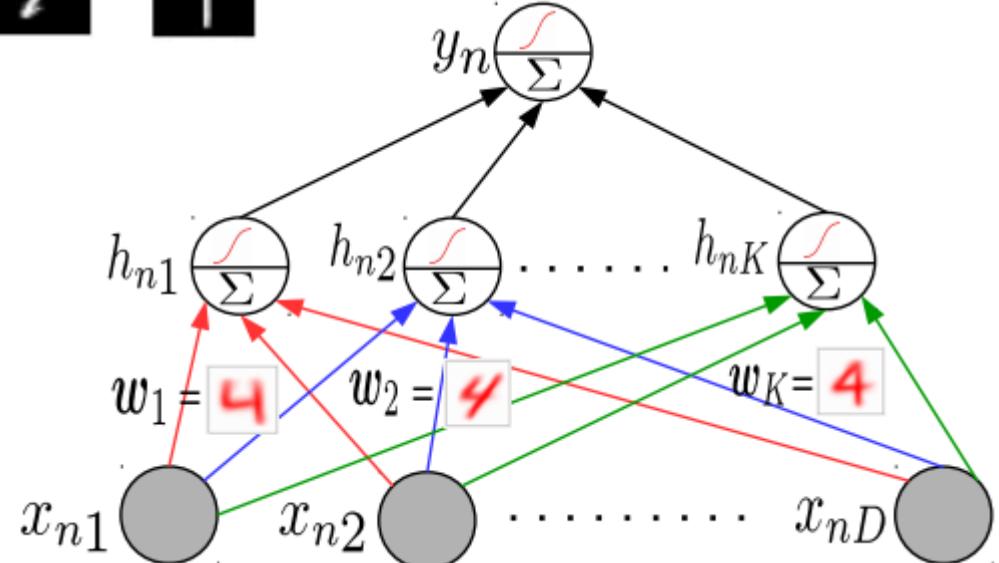
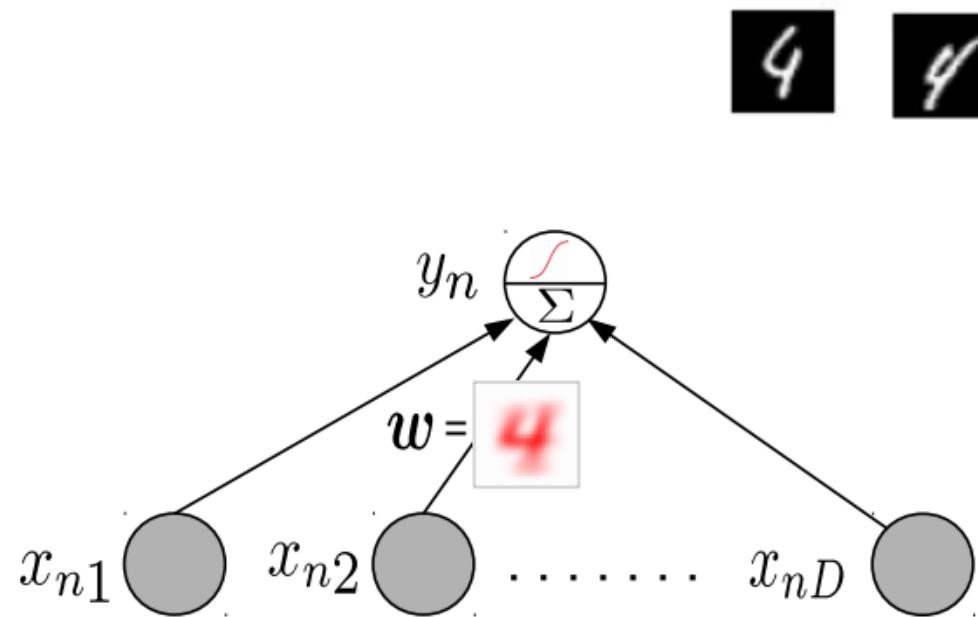
# Feature Learned by a Neural Network

- Node values in each hidden layer tell us how much a “learned” feature is active in  $\mathbf{x}_n$
- Hidden layer weights are like pattern/feature-detector/filter



# Why Neural Networks Work Better: Another View

- Linear models tend to only learn the “average” pattern
- Deep models can learn multiple patterns (each hidden node can learn one pattern)
  - Thus deep models can learn to capture more subtle variations that a simpler linear model



# Interconnection Among Nodes

Neural nets can differ in how nodes in one layer interconnect to the next layer (tonight)

Most General Case :- Fully connected.

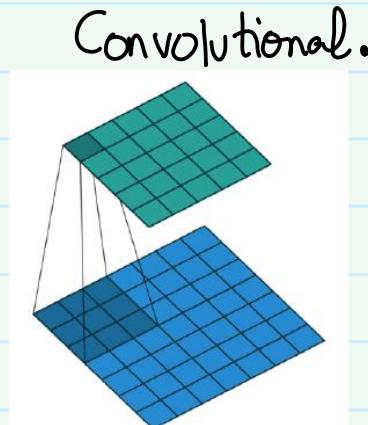
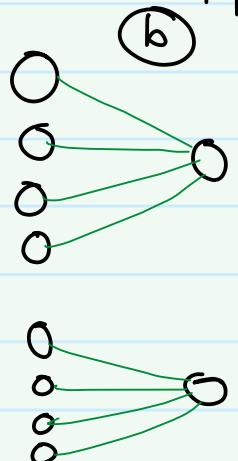
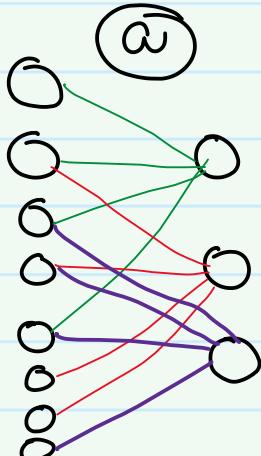
$$\text{# of Parameters} = (n_1 + 1)n_2$$

bias  
Left layer      Right layer nodes

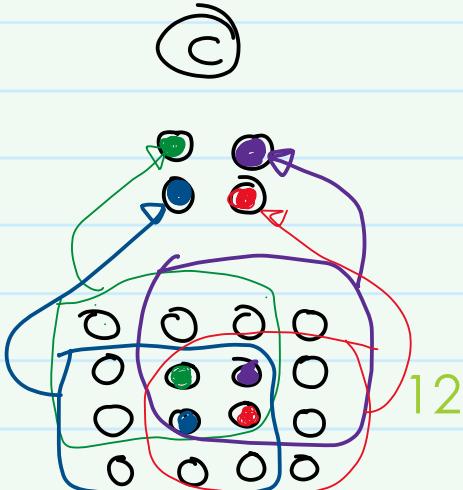
Other options :-

@ Random : For some 'm', pick m out of  $n_1$  nodes & make them input  $({}^n C_m)$  options

⑥ Pooled : Partition nodes of one layer into some # of clusters. In next layer (pooling layer) there is one node for each cluster & this node has all & only the member of its clusters as inputs.



Convolutional.



12

# Design Issues for Neural Nets

- Building a neural net to solve a given problem is partially art and partially science
- Before we begin to train a net by finding the weights on the inputs that serve our goals, we have to make a number of design decisions.
  - How many hidden layers should we use?
  - How many nodes will there be in each of the hidden layers?
  - In what manner will we inter-connect the outputs of one layer to the inputs of the next layer?
  - What activation function to use at each node?
  - What cost function should we minimize to express what weights are the best?
  - How to compute the outputs of each node as a function of the inputs?
- What algorithms shall we use to exploit the training examples in order to optimize the weights?

# Dense Feed-forward Networks

\* Ability to learn weights from training data.

Linear Algebra Notation.

Example problem :- Consecutive 1's in 4 bit vectors.

$$\text{Input, } x = [x_1 \ x_2 \ x_3 \ x_4]$$

$$\omega_{11} = [\omega_{111} \ \omega_{112} \ \omega_{113} \ \omega_{114}]$$

$$\omega_{12} = [\omega_{121} \ \omega_{122} \ \omega_{123} \ \omega_{124}]$$

$$\text{hidden node vector } h = [h_1 \ h_2]$$

$$h_1 = \text{sgn}(\omega_{11}^T x + b_1)$$

$$h_2 = \text{sgn}(\omega_{12}^T x + b_2)$$

$$\text{i.e., } h_i = \text{sgn}(\omega_i^T x + b_i) \text{ for } i=1,2$$

By organizing  $\omega_{11}$  &  $\omega_{12}$  in a  $2 \times 4$  matrix  $W$ , where  $i$ th row is  $\omega_i^T$

$$h = \text{sgn}(W x + b)$$

↑ works element-wise on vectors

$$y = \text{sgn}(w_2 h + c)$$

# Information Flow

Each Layer has its own additional matrix of weights & vector of biases  
+ Activation functions.

Perceptions used sign function, but other choices exist.

Since edges are directed forward. Information flows from left to right.

$\xrightarrow{\hspace{1cm}}$   
Feed forward networks. (No cycles)

Suppose there are  $l$  hidden layers & additional output layer  $l+1$

Let weight matrix for  $i^{\text{th}}$  layer be  $w_i$  & bias vector be  $b_i$ .

$w_1, w_2, \dots, w_{l+1}$  &  $b_1, b_2, \dots, b_{l+1}$   
These are the parameters of the model.

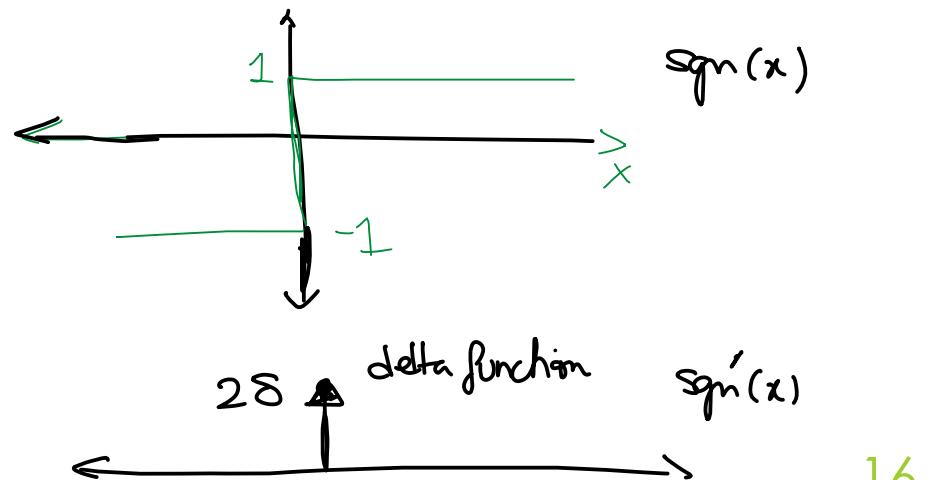
# Activation Functions

- $y = F_{l+1} (F_l(F_{l-1} \dots (F_2(F_1(x) + b_1) + b_2) \dots + b_{l-1}) + b_l) + b_{l+1}$
- Since we typically use gradient descent to solve for optimal value of parameters, we look for activation functions with the following properties
  1. The function is continuous and differentiable everywhere (or almost everywhere)
  2. The derivative of the function doesn't saturate (i.e. become very small, tending to zero) over its expected input range. Very small derivatives tend to stall the learning process
  3. The derivative of the function doesn't explode (i.e. become very large, tending to infinity), since this would lead to issues of numerical instability

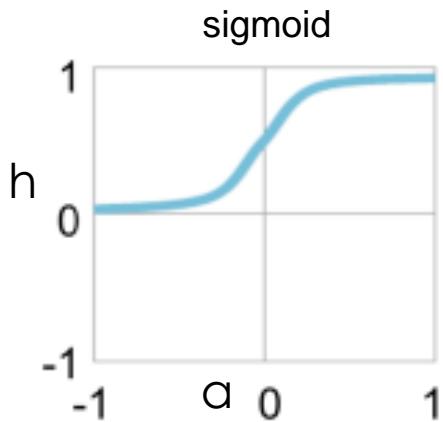
The sign function doesn't satisfy conditions 2, 3.

It's derivatives explode at 0 and is 0 everywhere else.

Let's therefore consider other activation functions.

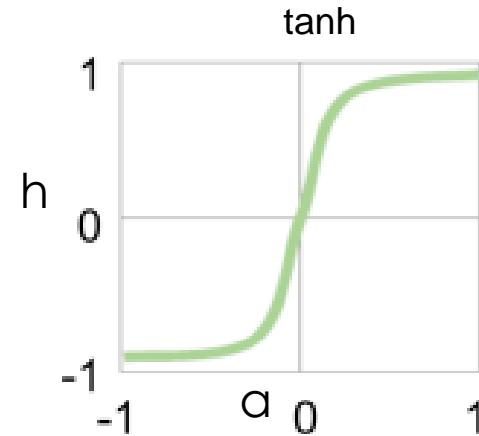


# Activation Functions: Some Common Choices



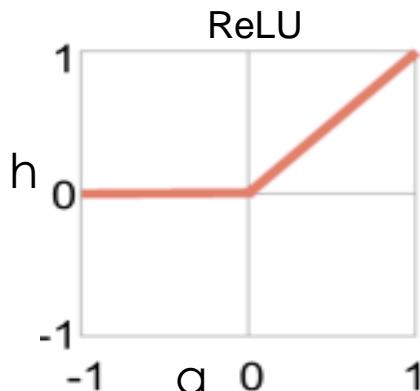
Sigmoid:  $h = \sigma(a) = \frac{1}{1+\exp(-a)}$

For sigmoid as well as tanh, gradients saturate (become close to zero as the function tends to its extreme values)



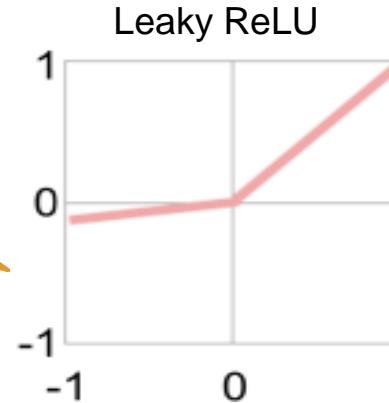
tanh (tan hyperbolic):  $h = \frac{\exp(a) - \exp(-a)}{\exp(a) + \exp(-a)} = 2\sigma(2a) - 1$

Preferred more than sigmoid. Helps keep the mean of the next layer's inputs close to zero (with sigmoid, it is close to 0.5)



ReLU (Rectified Linear Unit):  $h = \max(0, a)$

Helps fix the dead neuron problem of ReLU when  $a$  is a negative number



Leaky ReLU:  $h = \max(\beta a, a)$   
where  $\beta$  is a small positive number

ReLU and Leaky ReLU are among the most popular ones

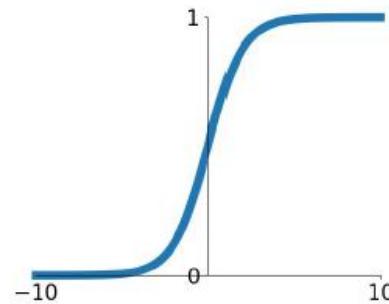
Without nonlinear activation, a deep neural network is equivalent to a linear model no matter how many layers we use



# Commonly Used Activation Functions

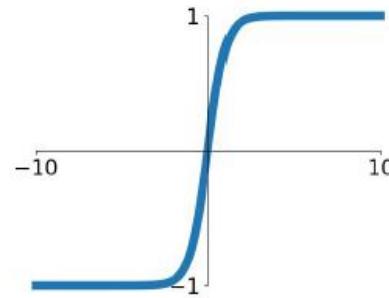
## Sigmoid

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



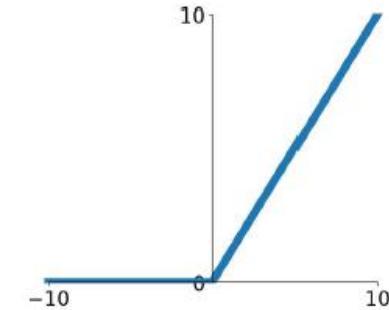
## tanh

$$\tanh(x) = \frac{e^{2x}-1}{e^{2x}+1}$$



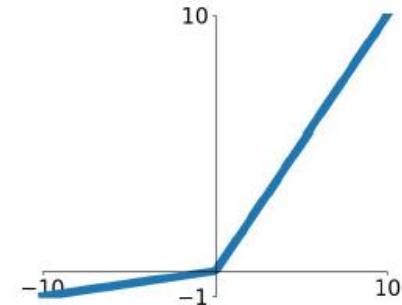
## ReLU

$$\max(0, x)$$



## Leaky ReLU

$$\max(0.1x, x)$$

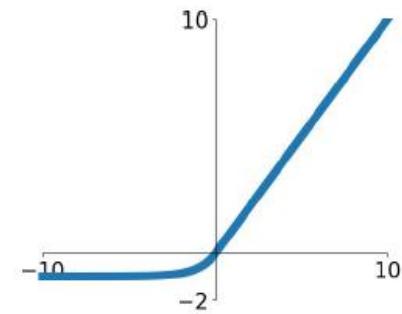


## Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

## ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



# Softmax

1 Operates on the entire vector.

If  $x = [x_1, x_2, \dots, x_n]$ , then its soft-max  $M(x) = [M(x_1), M(x_2) \dots M(x_n)]$

where

$$M(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

Pushes the largest component of the vector towards 1.  
All others towards 0.

Can be interpreted as a probability distribution.

$$\sum_i M(x_i) = 1$$

2 Common application : Output layer for a classification problem.

Output vector has a component for each class.

3 Works well with cross-entropy loss to handle saturation {very small gradients}

Accurate Calculation trick :-  $M(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}} = \frac{e^{x_i - c}}{\sum_j e^{x_j - c}}$  for any  $c$

Choose  $c = \max_j x_j$  so that  $x_j - c \leq 0$  for all  $j$ . Thus  $\forall j \quad 0 \leq e^{x_j - c} \leq 1$  19

# Cross Entropy Loss

Consider a multi-class classification problem with target classes  $C_1, C_2, C_3 \dots C_n$ . Suppose each point in training set is of the form  $(x, p)$  where  $x$  is input &  $p = [p_1 \ p_2 \ p_3 \ \dots \ p_n]$  is output.  $p_i$  is the probability that input  $x$  belongs to class  $C_i$  with  $\sum p_i = 1$ .

If we are certain that an input belongs to a particular class  $C_i$  then  $p_i = 1$  &  $p_j = 0 \ \forall j \neq i$ .

We can design our neural-net to produce output  $q = [q_1 \ q_2 \ \dots \ q_n]$ , using softmax for eg.

Hence both input & output are probability distributions.

Recall that entropy =  $\sum_i -p_i \log p_i = H(p)$

$$\text{Cross entropy} = -\sum_i p_i \underbrace{\log q_i}_{\substack{\text{Actual probability} \\ (\text{ground truth})}} = H(p, q)$$

Estimated probability

In information theory  
it is used to quantify  
Sub-optimal codes.

In general

$$H(p, q) \geq H(p)$$

KL Divergence

$$D(p||q) = H(p, q) - H(p) \\ = \sum_i p_i \log \frac{p_i}{q_i}$$

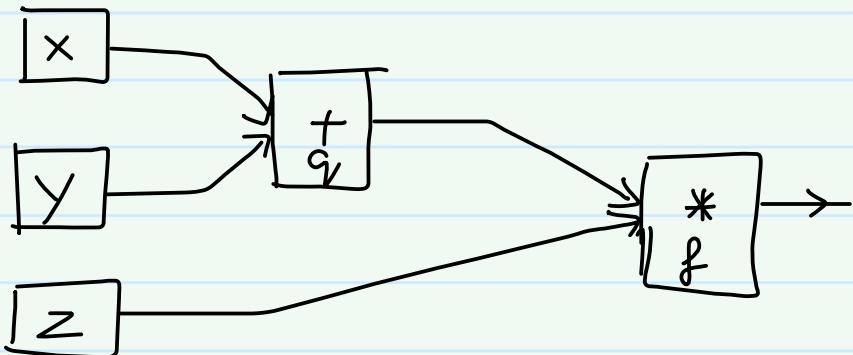
# Compute Graph

Find good values for parameters of the model. [weights & thresholds]

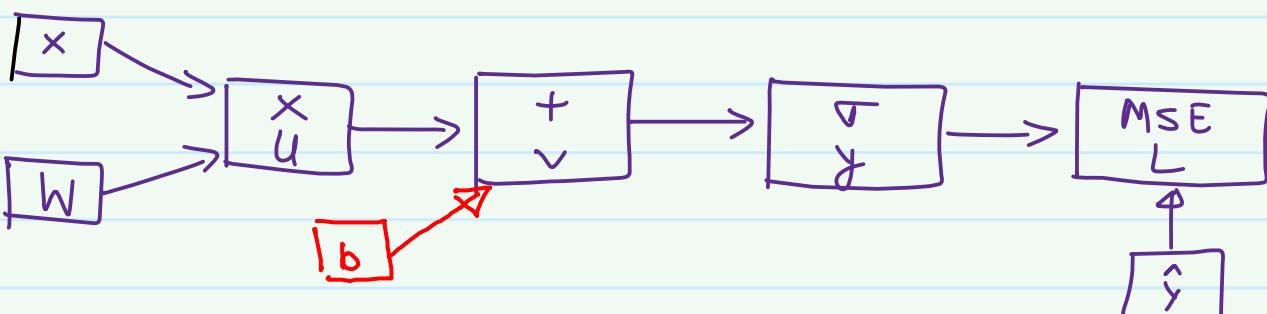
Compute Graph :- Directed graph, acyclic [DAG]

Each node has an operand or an operator (+, -,  $\sigma$ , MSE)

When a node has both an operand & operator, operator is written above the operand



$$f = (x + y)z$$



$$\begin{aligned} y &= \sigma(wx+b) \\ u &= wx \\ v &= u+b \\ y &= \sigma(v) \\ L &= \text{MSE}(y, \hat{y}) \end{aligned}$$

# Example of Backprop

$$f(x, y, z) = (x+y)z$$

$$q = x + y$$

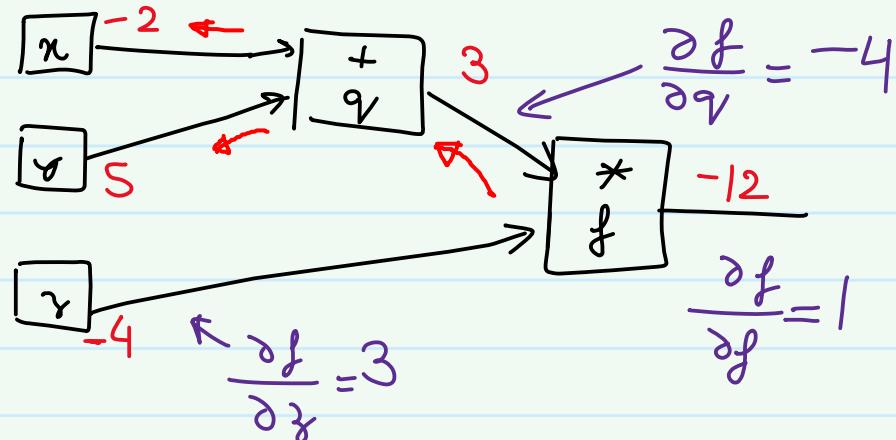
$$\frac{\partial q}{\partial x} = 1$$

$$\frac{\partial q}{\partial y} = 1$$

$$f = q \cdot z$$

$$\frac{\partial f}{\partial q} = z \quad \frac{\partial f}{\partial z} = q$$

Compute  $\frac{\partial f}{\partial x}$   $\frac{\partial f}{\partial y}$   $\frac{\partial f}{\partial z}$



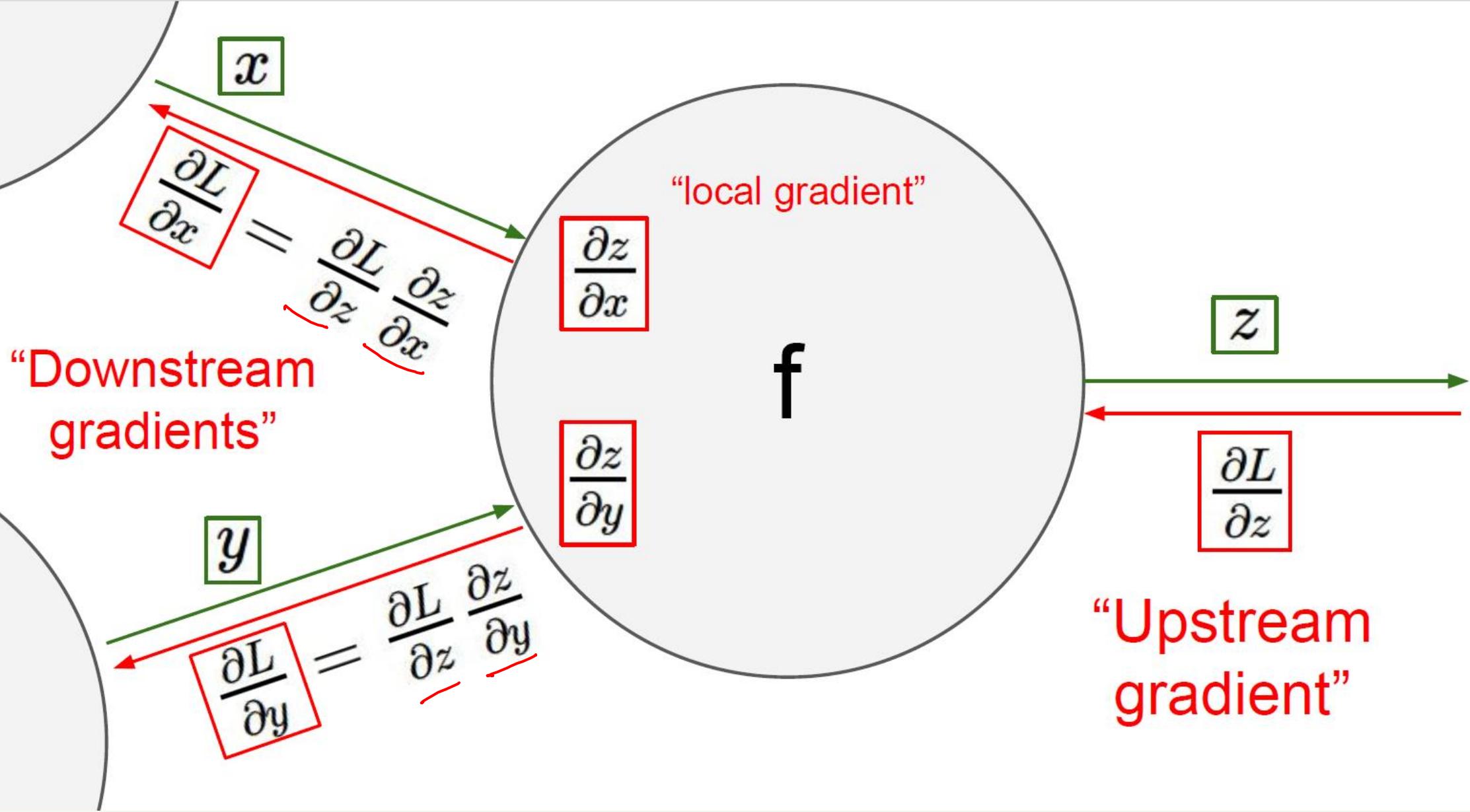
*Upstream gradient*  $\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \cdot \frac{\partial q}{\partial y}$  *Local gradient*

$$-4 \cdot 1 = -4$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \cdot \frac{\partial q}{\partial x}$$

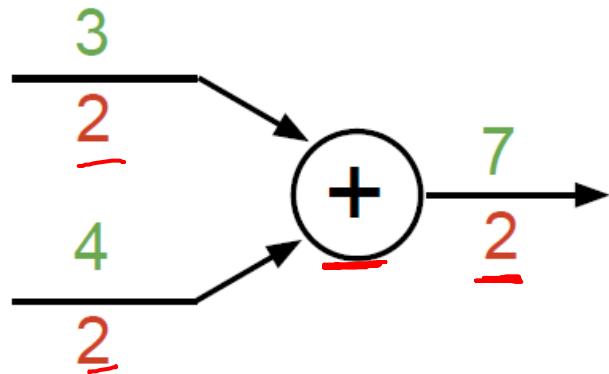
$$= -4 \cdot 1 = -4$$

# Upstream and Local Gradients

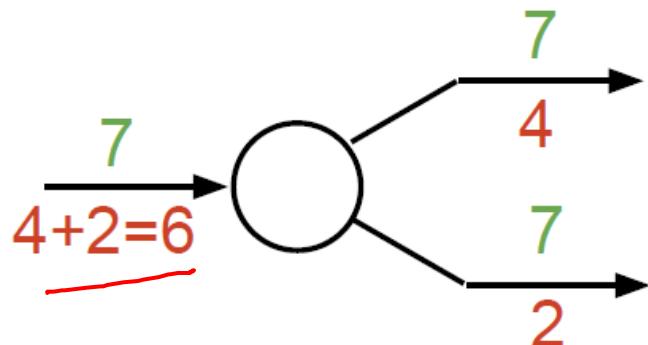


# Patterns in Gradient Flow

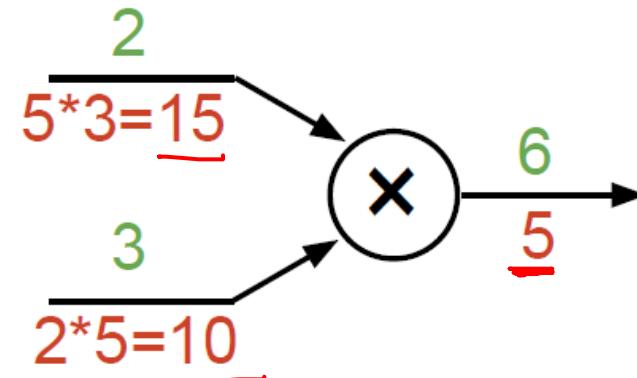
**add** gate: gradient distributor



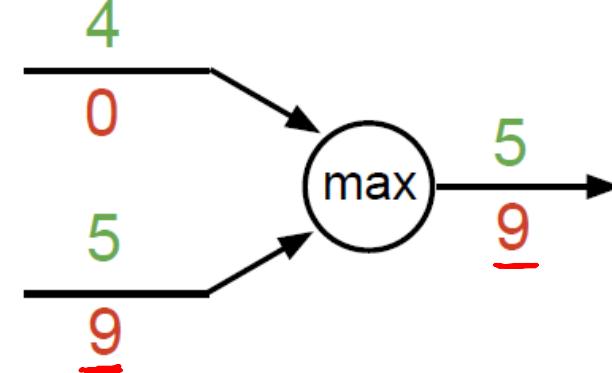
**copy** gate: gradient adder



**mul** gate: “swap multiplier”



**max** gate: gradient router



# Bonus Question 14.1 (5 marks)

\*  $f(\omega, x) = \frac{1}{1 + e^{-(\omega_0 x_0 + \omega_1 x_1 + \omega_2)}}$

Draw the compute graph for this function.

Assume  $\underline{\omega_0 = 2}$ ,  $\underline{\omega_1 = -3}$ ,  $\underline{x_0 = -1}$ ,  $x_1 = -2$ ,  $\underline{\omega_2 = -3}$  are input nodes

Now perform the forward & backprop to compute gradients w.r.t. each input  $\underline{x_i \text{ & } \omega_i}$

# Gradients, Jacobians, and Chain Rule

Goal :- Compute Gradient of loss function w.r.t. parameters of the network.

Recall :- If  $f: \mathbb{R}^N \rightarrow \mathbb{R}$  from real-valued vector to scalar  
 $x = [x_1, x_2, \dots, x_n] \quad \curvearrowright y = f(x)$

$$\nabla_x y = \left[ \frac{\partial y}{\partial x_1}, \frac{\partial y}{\partial x_2}, \dots, \frac{\partial y}{\partial x_n} \right]$$

Vector Valued :-  $f: \mathbb{R}^N \rightarrow \mathbb{R}^m$

Chain Rule

$$J = \begin{bmatrix} (\nabla f_1)^T \\ (\nabla f_2)^T \\ \vdots \\ (\nabla f_m)^T \end{bmatrix}$$

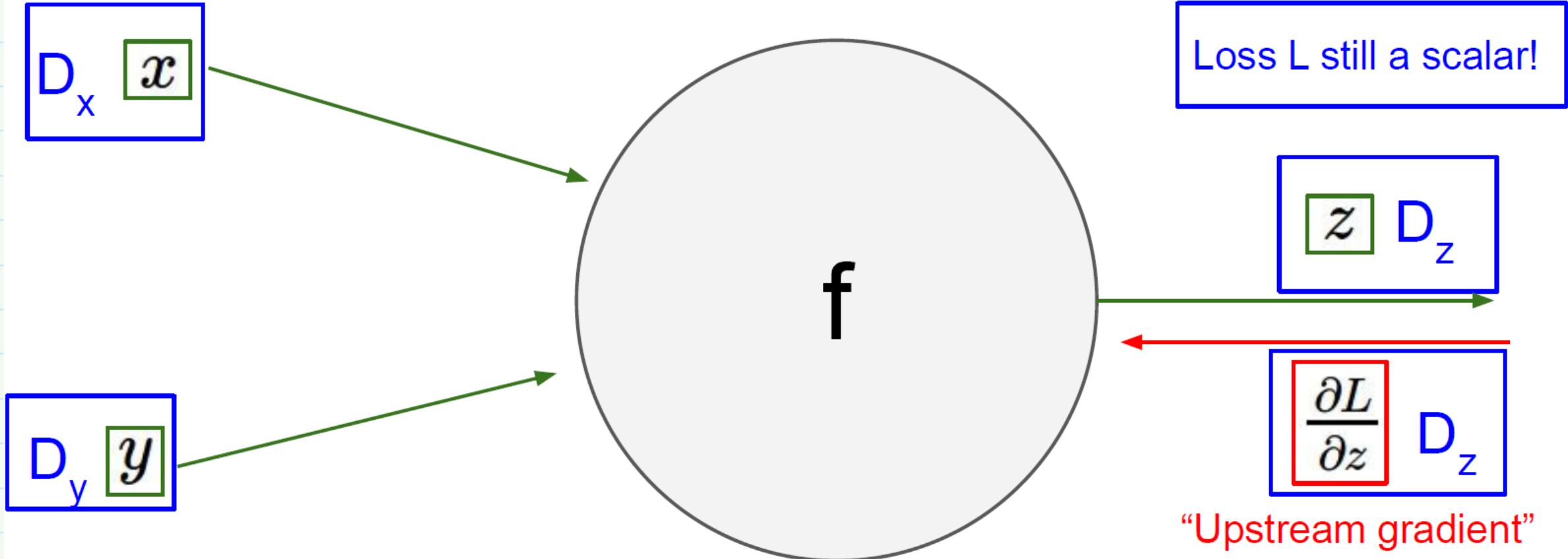
$$y = g(x), z = f(y) = f(g(x))$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

If  $z = f(u, v)$  where  $u = g(x)$  &  $v = h(x)$

$$\frac{dz}{dx} = \frac{\partial z}{\partial u} \cdot \frac{du}{dx} + \frac{\partial z}{\partial v} \cdot \frac{dv}{dx}$$

# Backprop with vectors



Loss  $L$  still a scalar!

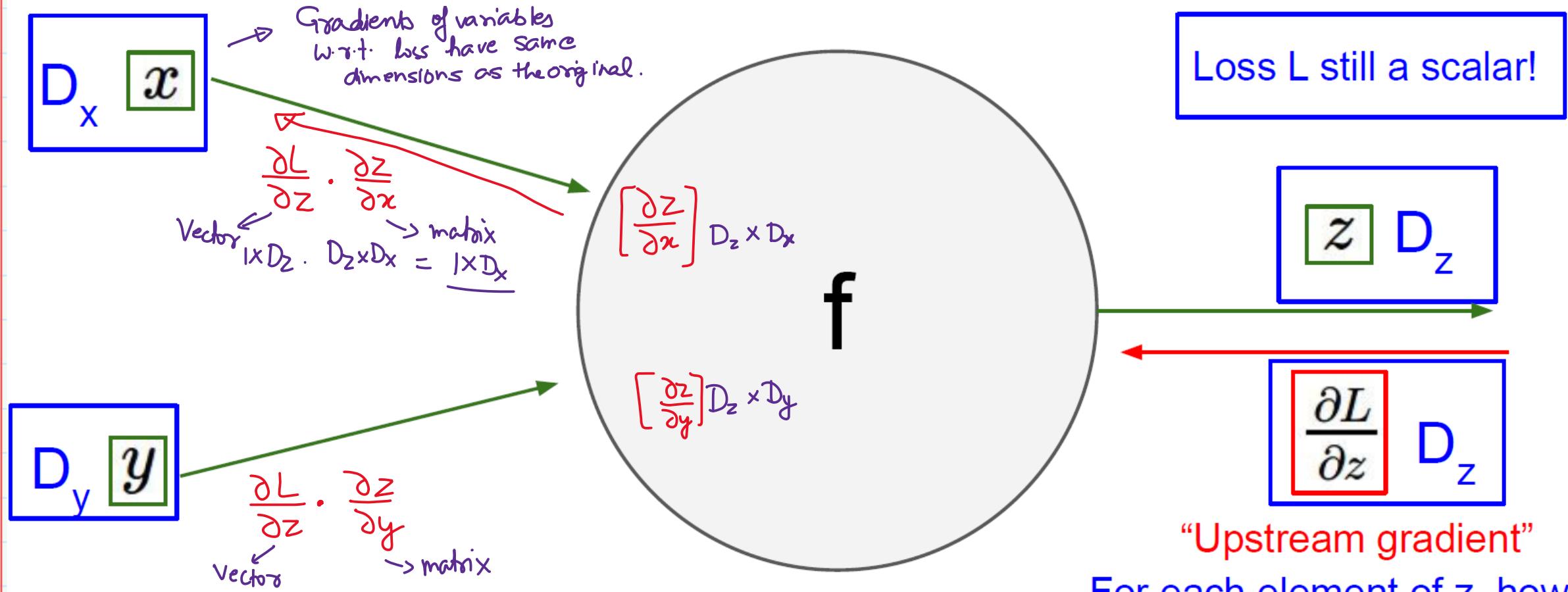


$$\frac{\partial L}{\partial z} \quad D_z$$

"Upstream gradient"

For each element of  $z$ , how much does it influence  $L$ ?

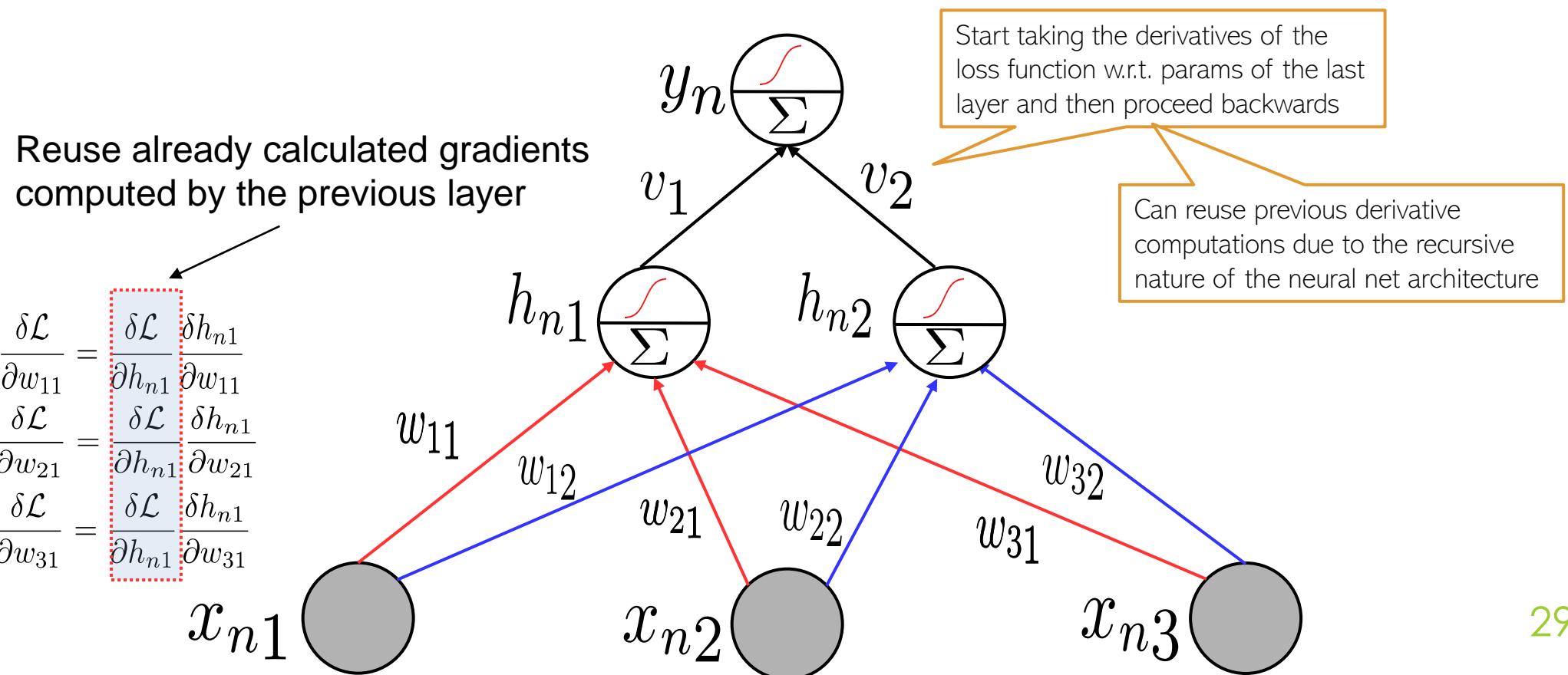
# Backprop with vectors



For each element of  $z$ , how much does it influence  $L$ ?

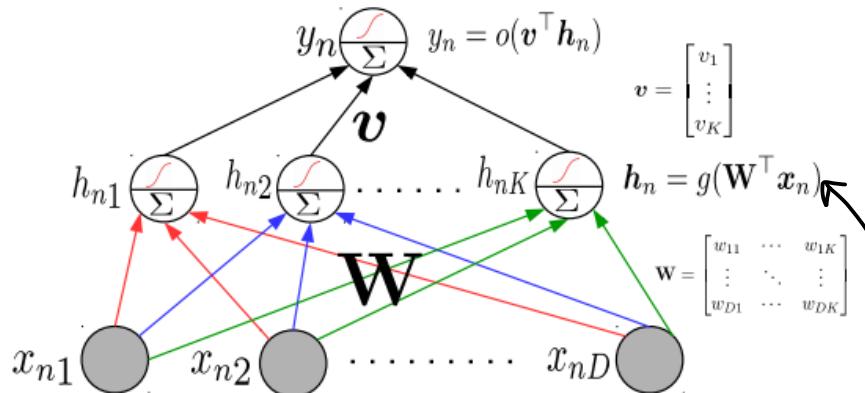
# Backpropagation with vectors

- Backpropagation = Gradient descent using chain rule of derivatives
- Chain rule of derivatives: Example, if  $y = f_1(x)$  and  $x = f_2(z)$  then  $\frac{\partial y}{\partial z} = \frac{\partial y}{\partial x} \frac{\partial x}{\partial z}$



# Backpropagation through an example

Consider a single hidden layer MLP



Assuming regression ( $o = \text{identity}$ ),  
the loss function for this model

$$\begin{aligned}\mathcal{L} &= \frac{1}{2} \sum_{n=1}^N (y_n - v^T h_n)^2 \\ &= \frac{1}{2} \sum_{n=1}^N \left( y_n - \sum_{k=1}^K v_k h_{nk} \right)^2 \\ &= \frac{1}{2} \sum_{n=1}^N \left( y_n - \sum_{k=1}^K v_k g(w_k^T x_n) \right)^2\end{aligned}$$

- To use gradient methods for  $\mathbf{W}, \mathbf{v}$ , we need gradients.
- Gradient of  $\mathcal{L}$  w.r.t.  $\mathbf{v}$  is straightforward

$$\frac{\partial \mathcal{L}}{\partial v_k} = - \sum_{n=1}^N \left( y_n - \sum_{k=1}^K v_k g(w_k^T x_n) \right) h_{nk} = \sum_{n=1}^N e_n h_{nk}$$

*error*

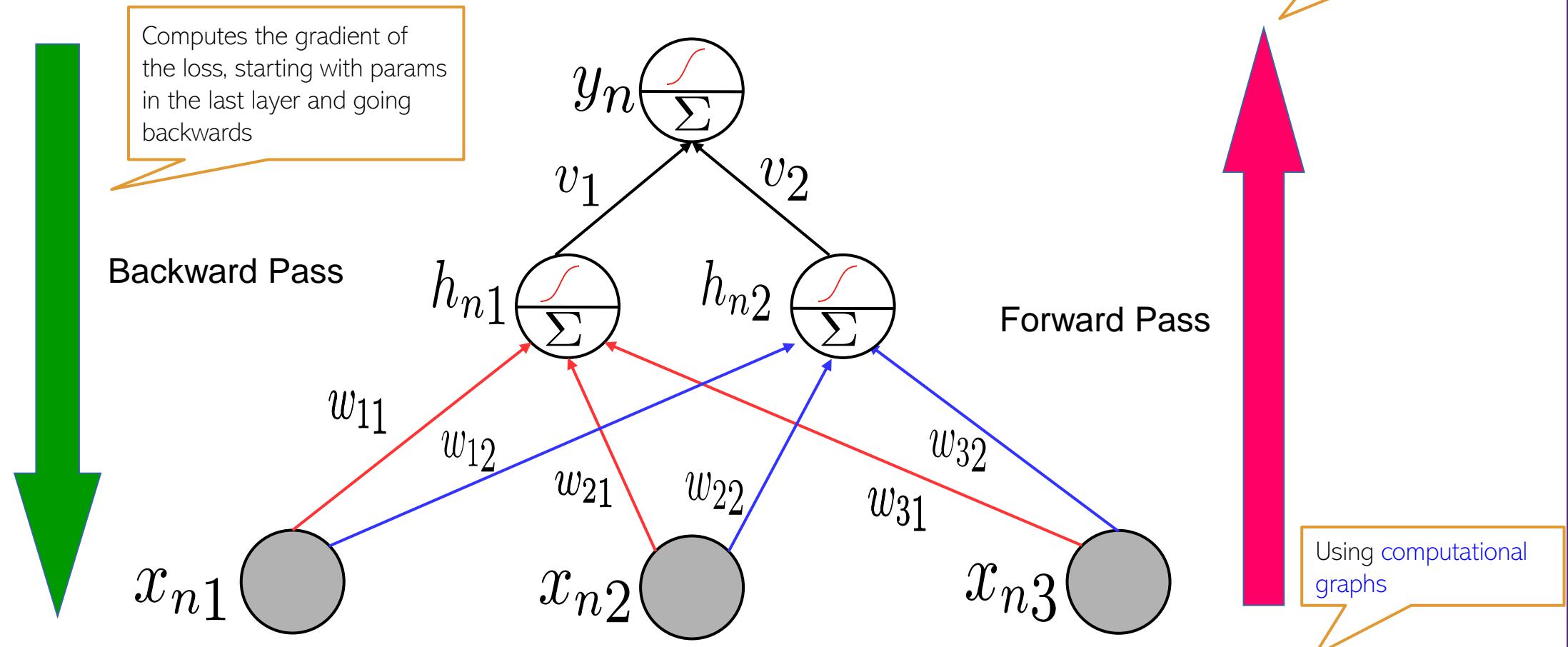
- Gradient of  $\mathcal{L}$  w.r.t.  $\mathbf{W}$  requires chain rule

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial w_{dk}} &= \sum_{n=1}^N \frac{\partial \mathcal{L}}{\partial h_{nk}} \frac{\partial h_{nk}}{\partial w_{dk}} \\ \frac{\partial \mathcal{L}}{\partial h_{nk}} &= -(y_n - \sum_{k=1}^K v_k g(w_k^T x_n)) v_k = -e_n v_k \\ \frac{\partial h_{nk}}{\partial w_{dk}} &= g'(w_k^T x_n) x_{nd} \quad (\text{note: } h_{nk} = g(w_k^T x_n))\end{aligned}$$

- Forward prop computes errors  $e_n$  using current  $\mathbf{W}, \mathbf{v}$ .  
Backprop updates NN params  $\mathbf{W}, \mathbf{v}$  using grad methods
- Backprop caches many of the calculations for reuse

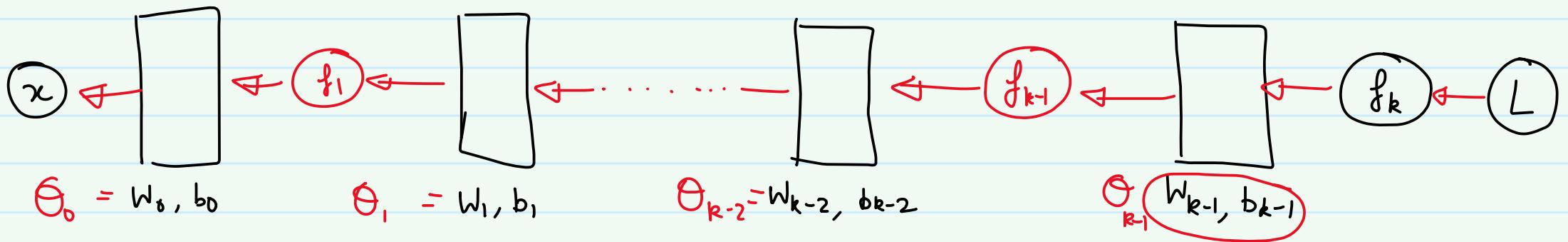
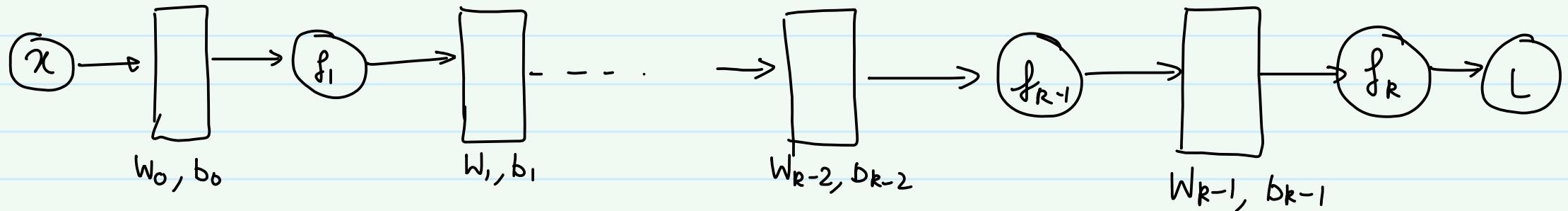
# Backpropagation

- Backprop iterates between a forward pass and a backward pass



Software frameworks such as Tensorflow and PyTorch support this already so you don't need to implement it by hand (so no worries of computing derivatives etc)

# Gradients in a deep neural net

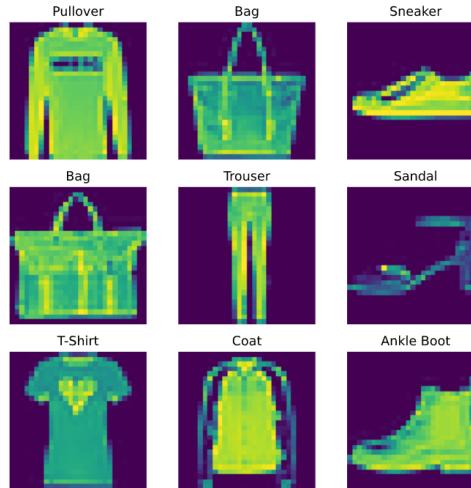
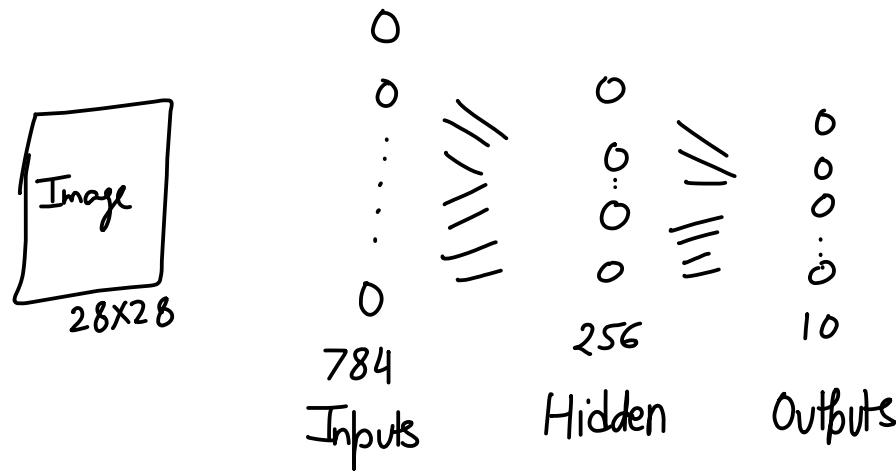


$$\frac{\partial L}{\partial \theta_{k-1}} = \frac{\partial L}{\partial f_k} \cdot \frac{\partial f_k}{\partial \theta_{k-1}}$$

$$\frac{\partial L}{\partial \theta_{k-2}} = \frac{\partial L}{\partial f_k} \cdot \frac{\partial f_k}{\partial f_{k-1}} \cdot \frac{\partial f_{k-1}}{\partial \theta_{k-2}}$$

$$\frac{\partial L}{\partial \theta_i} = \frac{\partial L}{\partial f_k} \cdot \frac{\partial f_k}{\partial f_{k-1}} \cdots \frac{\partial f_{i+2}}{\partial f_{i+1}} \cdot \frac{\partial f_{i+1}}{\partial \theta_i}$$

# Example of Feed-forward Neural Network



```
num_inputs, num_outputs, num_hiddens = 784, 10, 256
```

```
w1 = nn.Parameter(  
    torch.randn(num_inputs, num_hiddens, requires_grad=True) * 0.01)  
b1 = nn.Parameter(torch.zeros(num_hiddens, requires_grad=True))  
w2 = nn.Parameter(  
    torch.randn(num_hiddens, num_outputs, requires_grad=True) * 0.01)  
b2 = nn.Parameter(torch.zeros(num_outputs, requires_grad=True))  
  
params = [w1, b1, w2, b2]
```

```

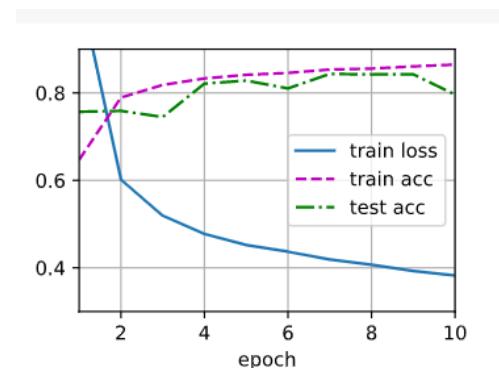
def relu(X):
    a = torch.zeros_like(X)
    return torch.max(X, a)

def net(X):
    X = X.reshape((-1, num_inputs))
    H = relu(X @ W1 + b1)  # Here '@' stands
    s for matrix multiplication
    return (H @ W2 + b2)

loss = nn.CrossEntropyLoss()

num_epochs, lr = 10, 0.1
updater = torch.optim.SGD(params, lr=lr)
d2l.train_ch3(net, train_iter, test_iter, loss, num_epochs, updater)

```



# Convolutional Neural Networks

Consider a fully connected net for processing  $1024 \times 1024$  [1 MP] images  $\times 3$  channels RGB

# of parameters for each output node =  $10^6 \times 3$

Suppose we have 320 nodes  $\sim 10^9$  (billion parameters)

Training set size usually  $10,000$  to  $10^5$ .

Result : Overfitting.

Solution :- Need regularization

- (a) Reduce connections
- (b) Share parameters across nodes.
- (c) Aggregate information

① Convolutional layers : Use local patterns

② Pooling : Aggregation over small spatial regions

Compute max of every  $2 \times 2$  region.

↳ Non overlapping regions

↳  $224 \times 224 \times 3 \Rightarrow 112 \times 112 \times 3$

5x5 Filter example

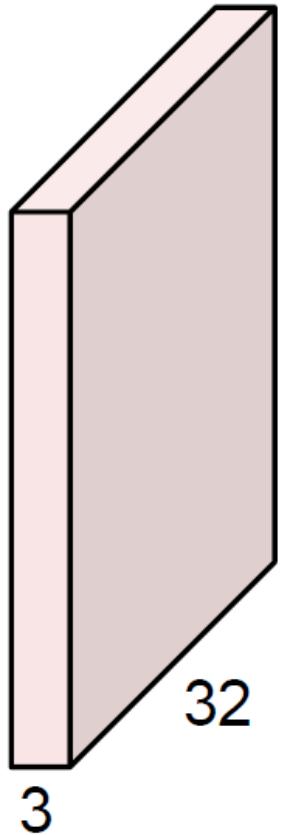
$$z_{ij} = \sum_{k=0}^4 \sum_{l=0}^4 x_{i+k, j+l} w_{kl} + b$$

↳ Can be passed through ReLU

# Create Trainable Filters

## Convolution Layer

32x32x3 image



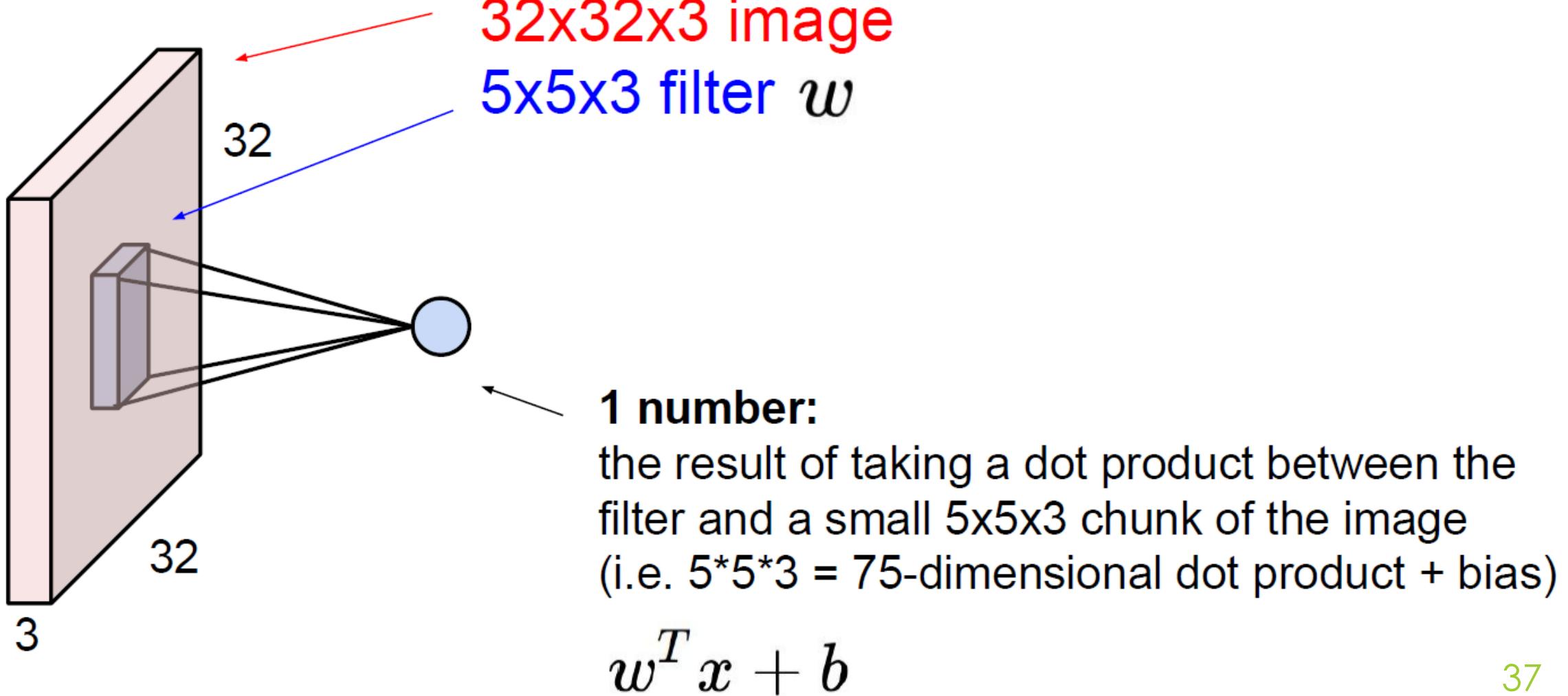
5x5x3 filter



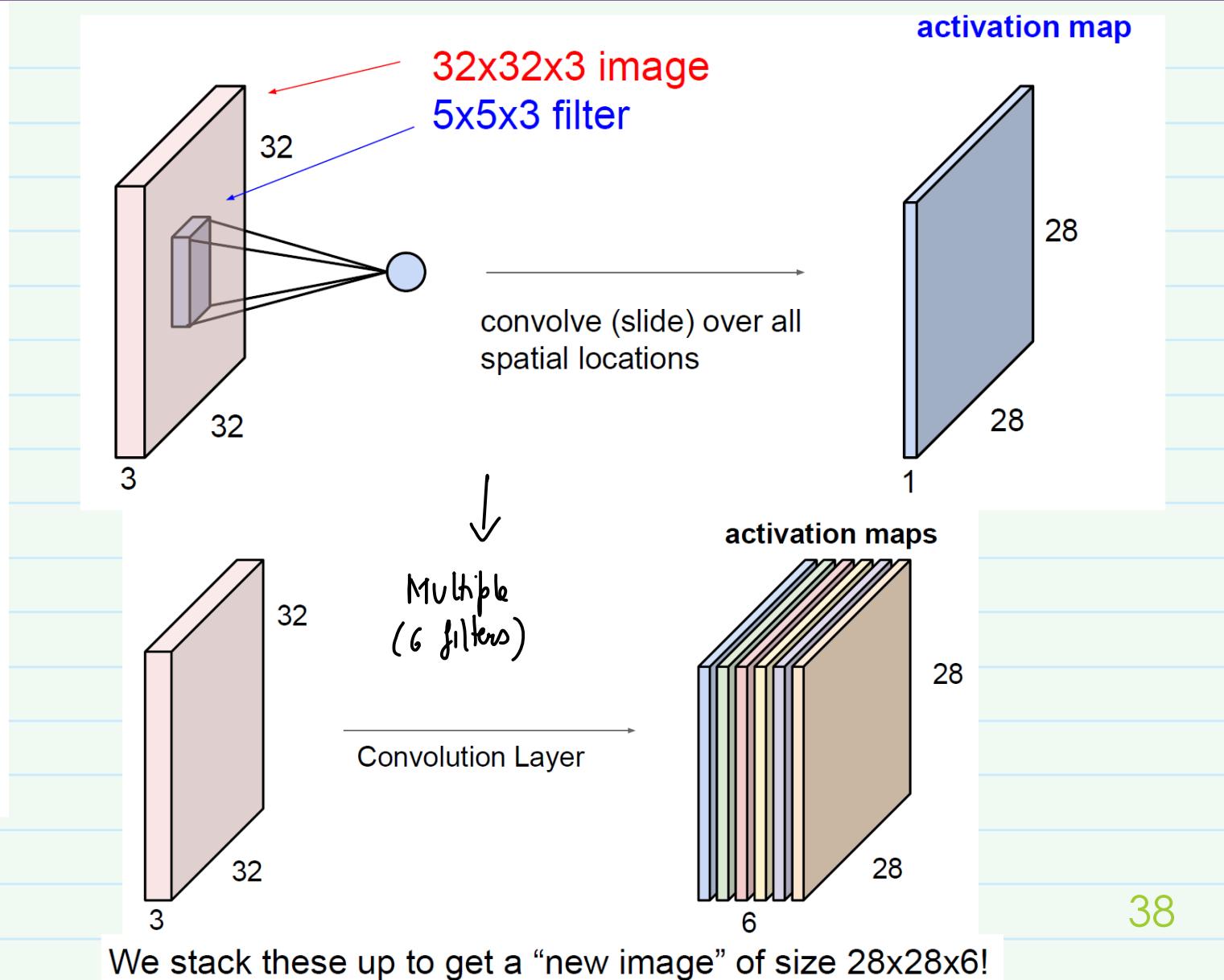
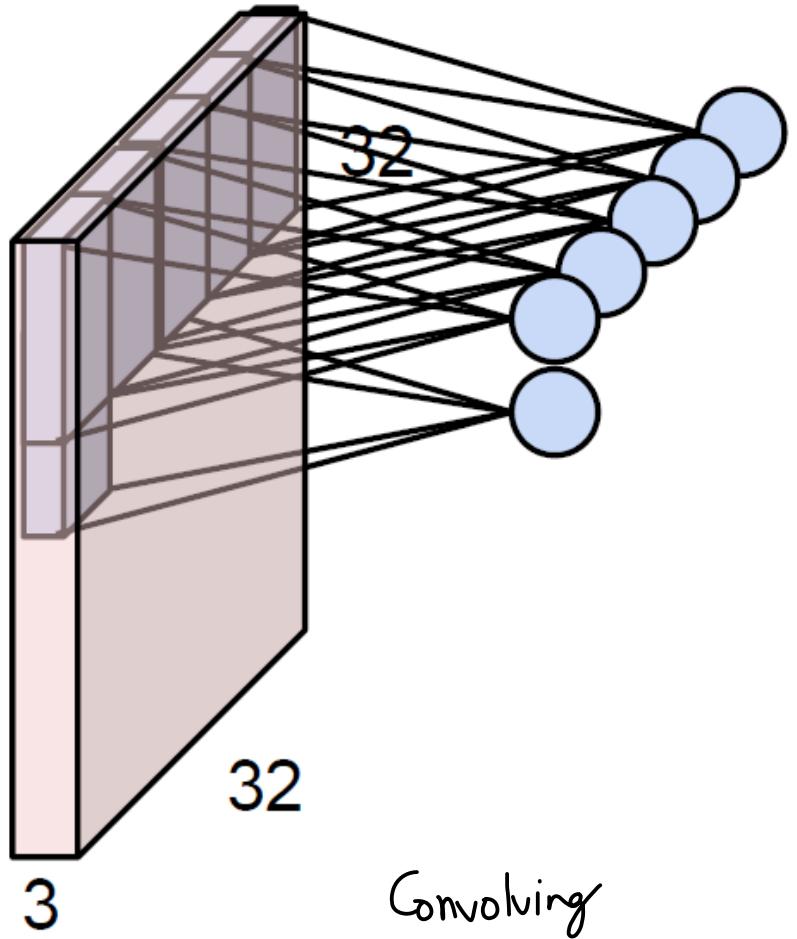
Filters always extend the full depth of the input volume

**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

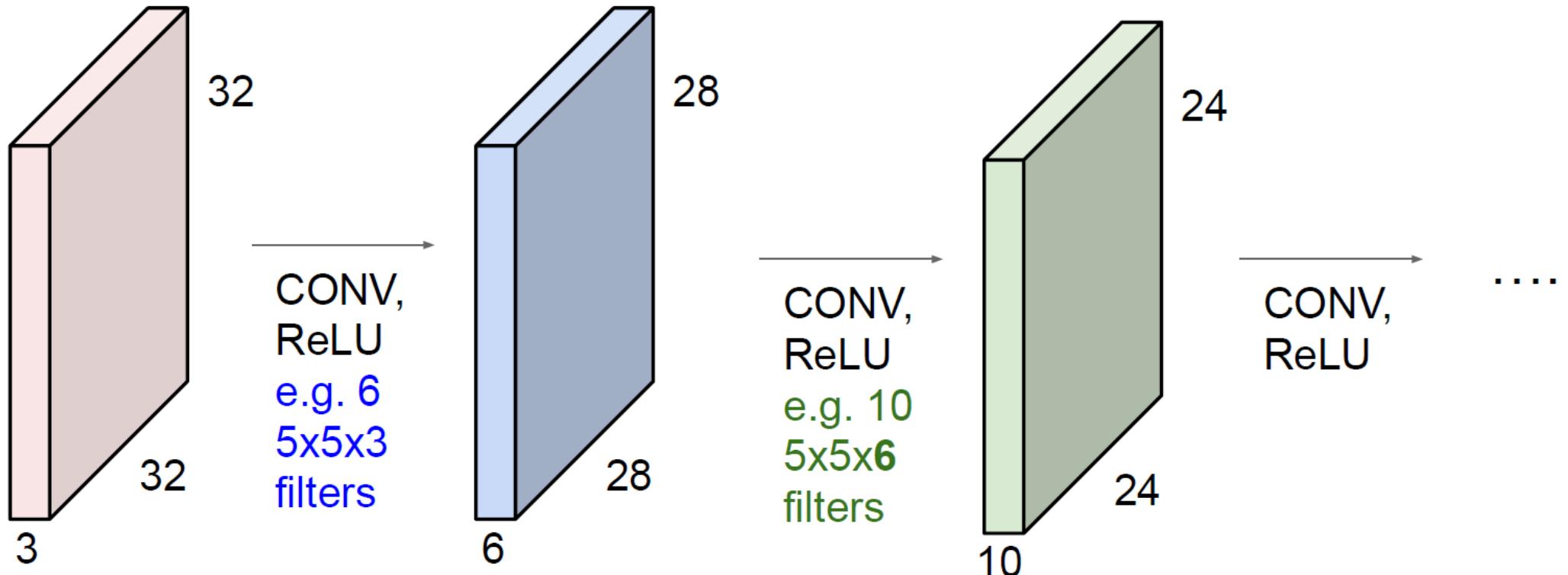
# Take Dot Products with the Filter



# Slide the filter all over the image to create maps

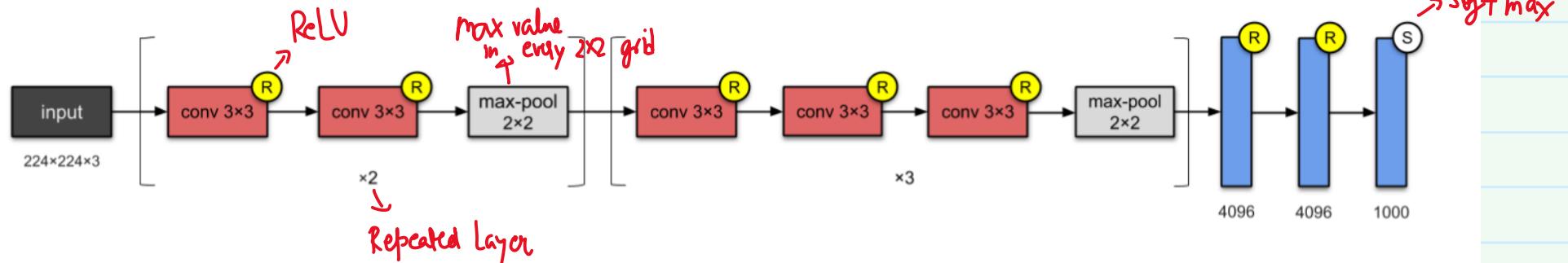
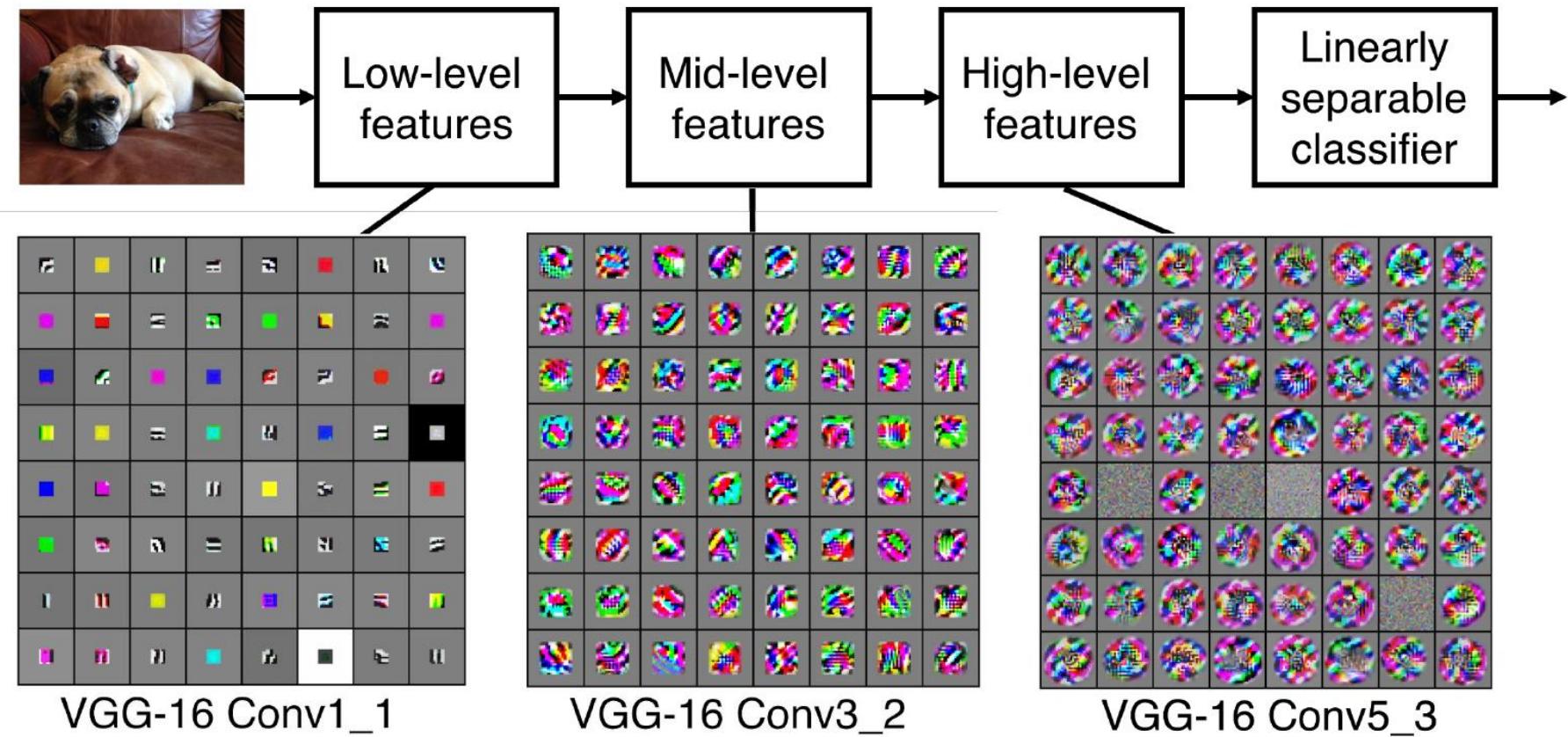


# Now we have a recipe!



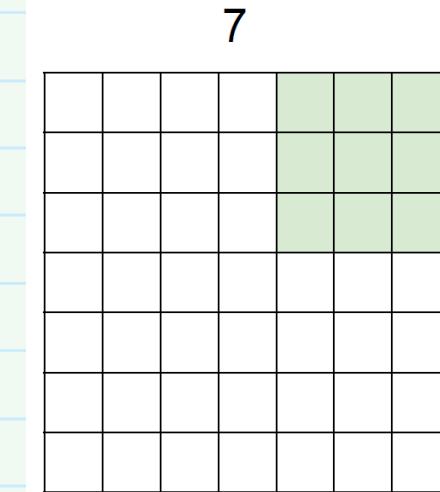
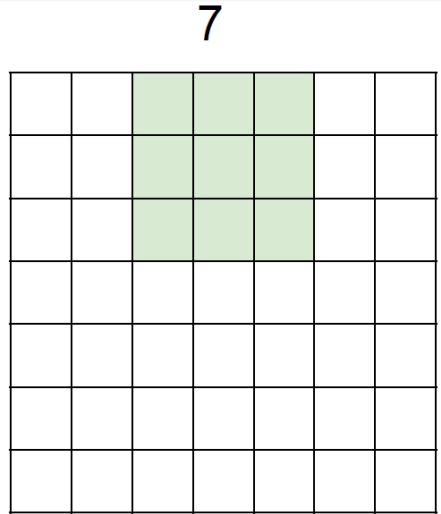
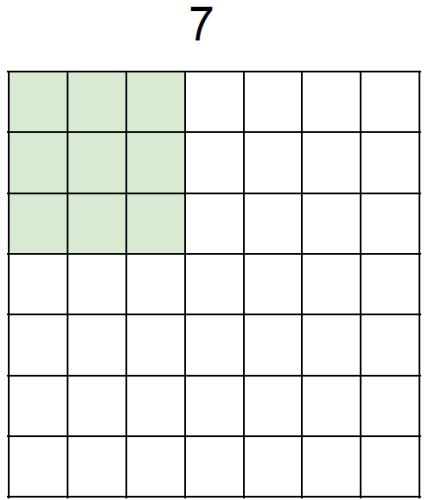
Convolutional Layers + Activation Functions.

# VGG-16 Design



which has 13 convolutional and 3 fully-connected layers, carrying with

# Using Strides



7

7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**  
**=> 3x3 output!**

Can we use 3x3 filter with stride 3?

No, doesn't fit!

Output size:  
**(N - F) / stride + 1**

e.g. N = 7, F = 3:

$$\text{stride 1} \Rightarrow (7 - 3)/1 + 1 = 5$$

$$\text{stride 2} \Rightarrow (7 - 3)/2 + 1 = 3$$

$$\text{stride 3} \Rightarrow (7 - 3)/3 + 1 = 2.33$$

# Zero Padding the border

0	0	0	0	0	0			
0								
0								
0								
0								

e.g. input 7x7

3x3 filter, applied with **stride 1**

**pad with 1 pixel border => what is the output?**

**7x7 output!**

(recall:)

$$(N + 2P - F) / \text{stride} + 1$$

in general, common to see CONV layers with stride 1, filters of size  $F \times F$ , and zero-padding with  $(F-1)/2$ . (will preserve size spatially)

e.g.  $F = 3 \Rightarrow$  zero pad with 1

$F = 5 \Rightarrow$  zero pad with 2

$F = 7 \Rightarrow$  zero pad with 3

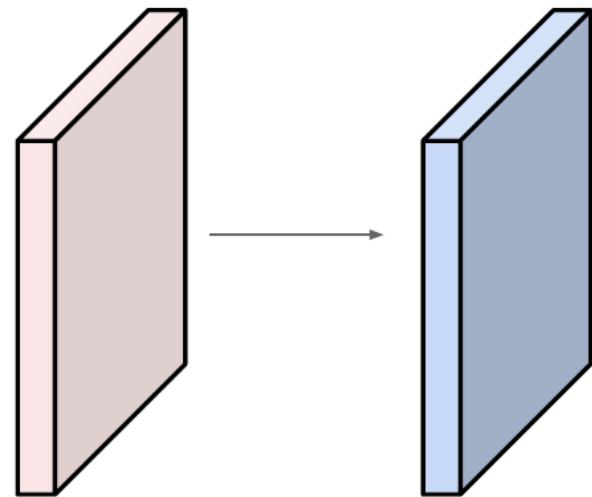
Shrinking the image size too fast  
isn't considered good.

# Calculating the number of parameters

Examples time:

Input volume: **32x32x3**

**10 5x5** filters with stride 1, pad 2



Number of parameters in this layer?

each filter has  $5*5*3 + 1 = 76$  params      (+1 for bias)

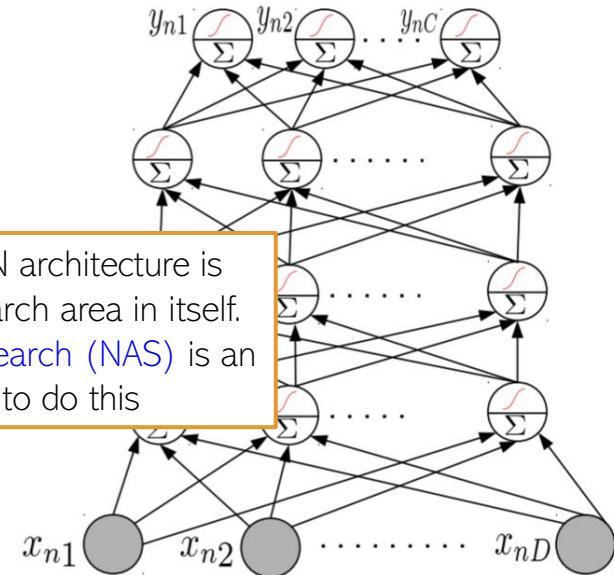
$$\Rightarrow 76 * 10 = 760$$



# Neural Nets: Some Design Aspects

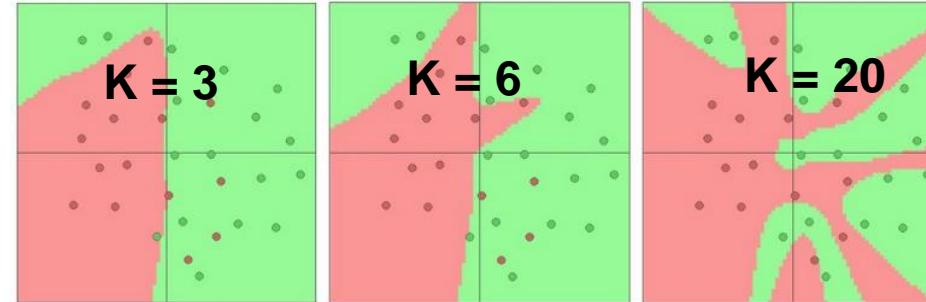
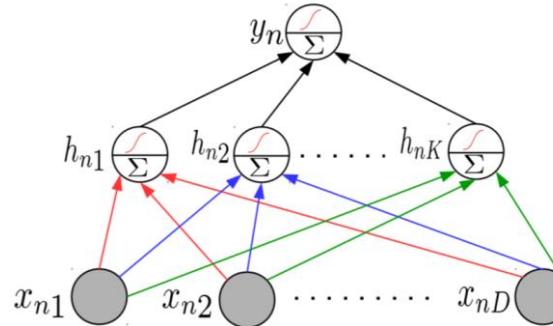
- Much of the magic lies in the hidden layers
- Hidden layers learn and detect good features
- Need to consider a few aspects
  - Number of hidden layers, number of units in each hidden layer
  - Why bother about many hidden layers and not use a single very wide hidden layer (recall Hornik's universal function approximator theorem)?
  - Complex networks (several, very wide hidden layers) or simpler networks (few, moderately wide hidden layers)?
  - Aren't deep neural network prone to overfitting (since they contain a huge number of parameters)?

Choosing the right NN architecture is important and a research area in itself. Neural Architecture Search (NAS) is an automated technique to do this



# Representational Power of Neural Nets

- Consider a single hidden layer neural net with  $K$  hidden nodes



- Recall that each hidden unit “adds” a function to the overall function
- Increasing  $K$  (number of hidden units) will result in a more complex function
- Very large  $K$  seems to overfit (see above fig). Should we instead prefer small  $K$ ?
- No! It is better to use large  $K$  and regularize well. Reason/justification:
  - Simple NN with small  $K$  will have a few local optima, some of which may be bad
  - Complex NN with large  $K$  will have many local optima, all equally good (theoretical results on this)
- We can also use multiple hidden layers (each sufficiently large) and regularize well<sup>46</sup>

# Preventing Overfitting in Neural Nets

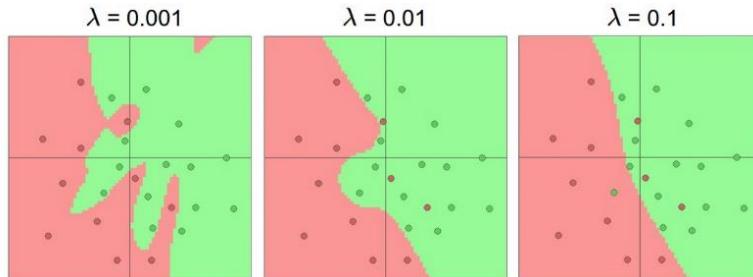
Various other tricks, such as weight sharing across different hidden units of the same layer (used in convolutional neural nets or CNN)

47



- Neural nets can overfit. Many ways to avoid overfitting, such as
  - Standard regularization on the weights, such as  $\ell_2$ ,  $\ell_1$ , etc ( $\ell_2$  reg. is also called weight decay)

Single Hidden Layer NN with K = 20 hidden units and L2 regularization



- Early stopping (traditionally used): Stop when validation error starts increasing
- Dropout: Randomly remove units (with some probability  $p \in (0,1)$ ) during training

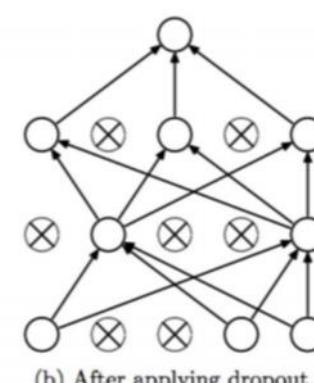
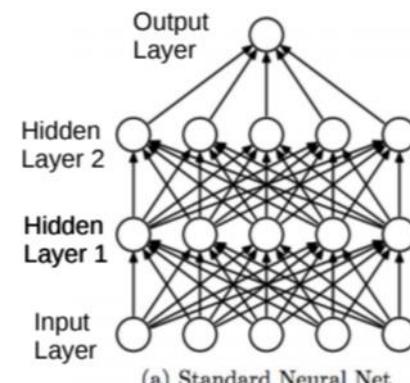
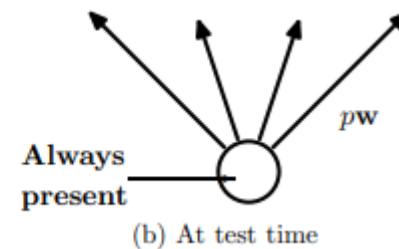
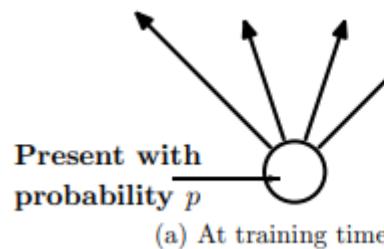
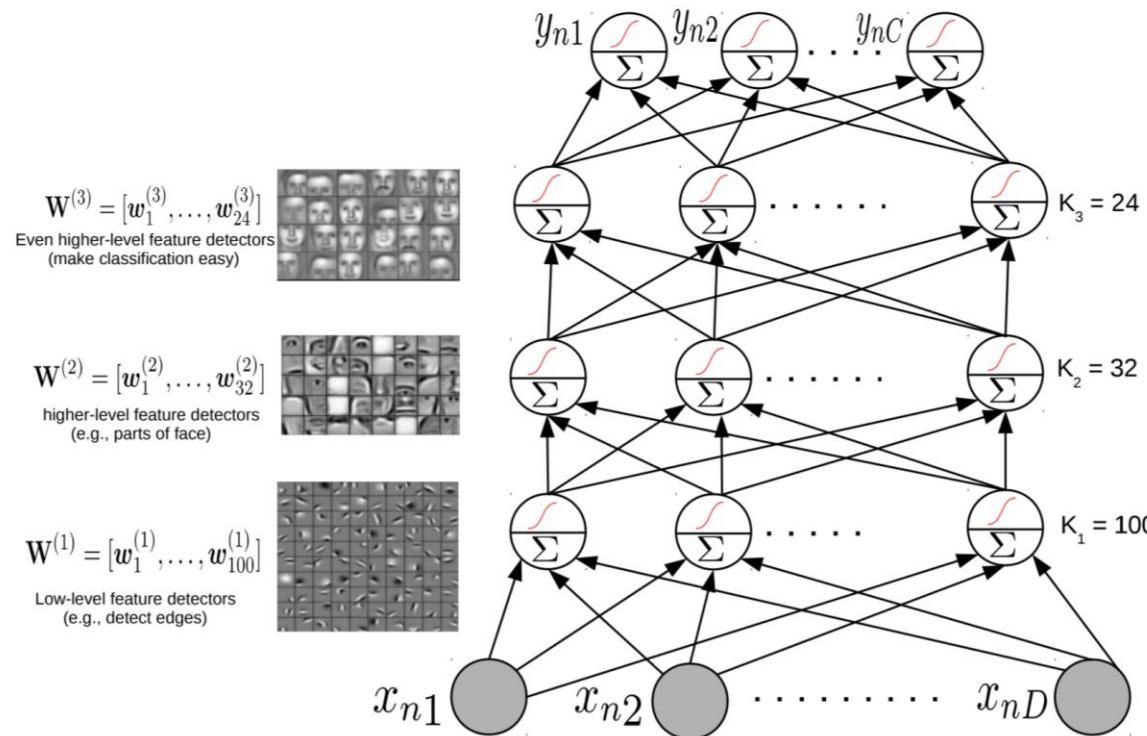


Fig courtesy: Dropout: A Simple Way to Prevent Neural Networks from Overfitting (Srivastava et al, 2014)

47

# Wide or Deep?

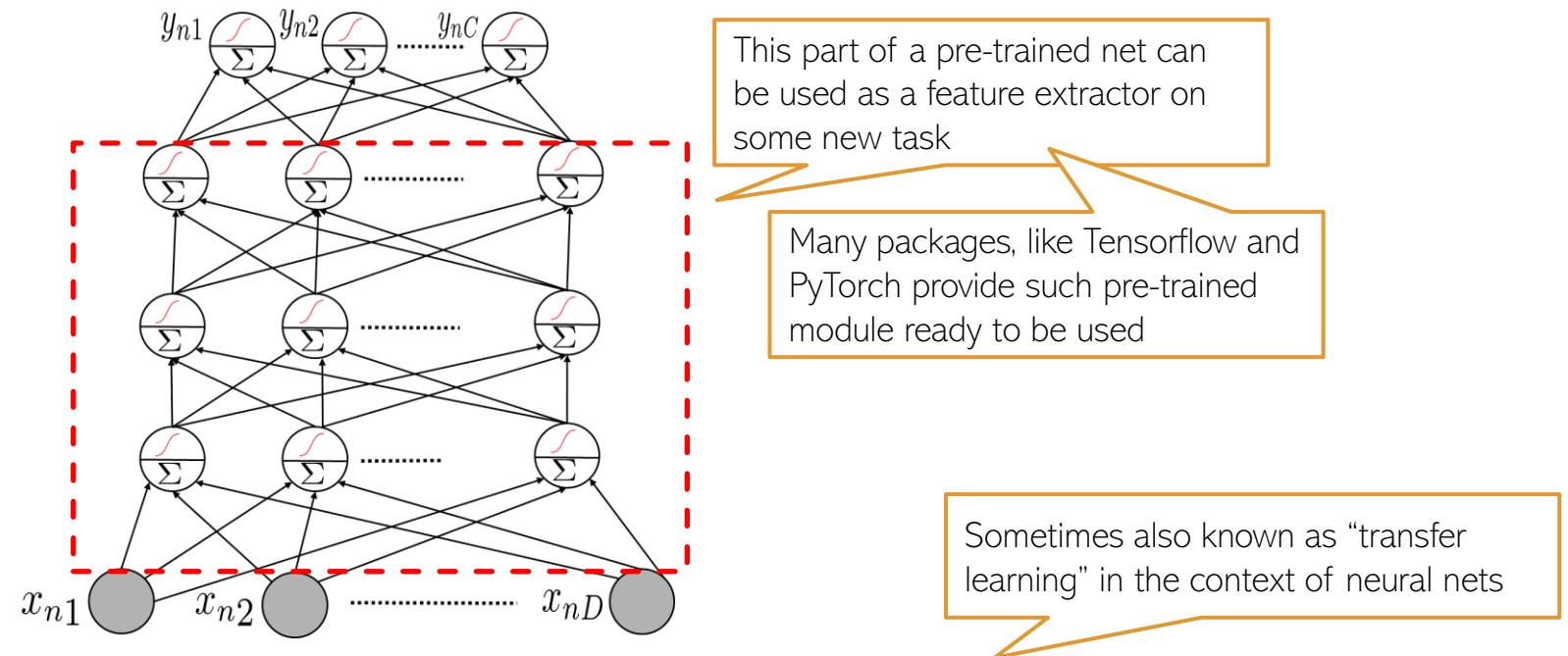
- While very wide single hidden layer can approx. any function, often we prefer many, less wide, hidden layers



- Higher layers help learn more directly useful/interpretable features (also useful for compressing data using a small number of features)

# Using a Pre-trained Network

- A deep NN already trained in some “generic” data can be useful for other tasks, e.g.,
  - Feature extraction: Use a pre-trained net, remove the output layer, and use the rest of the network as a feature extractor for a related dataset



- Fine-tuning: Use a pre-trained net, use its weights as initialization to train a deep net for a new but related task (useful when we don’t have much training data for the new task)

# Deep Neural Nets: Some Comments

- Highly effective in learning good feature rep. from data in an “end-to-end” manner
- The objective functions of these models are **highly non-convex**
  - But fast and robust non-convex opt algos exist for learning such deep networks
- Training these models is computationally very expensive
  - But GPUs can help to speed up many of the computations
- Also useful for unsupervised learning problems (will see some examples)
  - Autoencoders for dimensionality reduction
  - Deep generative models for generating data and (unsupervisedly) learning features – examples include generative adversarial networks (GAN) and variational auto-encoders (VAE)

# Bonus Question 2.1 Solutions

- $f_1 = \sin(x_1) \cos(x_2), \quad x \in \mathbb{R}^2$

$$\frac{\partial f_1}{\partial x_1} = \cos(x_1) \cos(x_2)$$

$$\frac{\partial f_1}{\partial x_2} = -\sin(x_1) \sin(x_2)$$

$$\implies J = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} \end{bmatrix} = [\cos(x_1) \cos(x_2) \quad -\sin(x_1) \sin(x_2)] \in \mathbb{R}^{1 \times 2}$$