

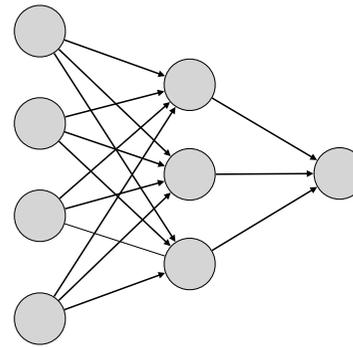
Neural Networks

Lecture 11

Termeh Shafie

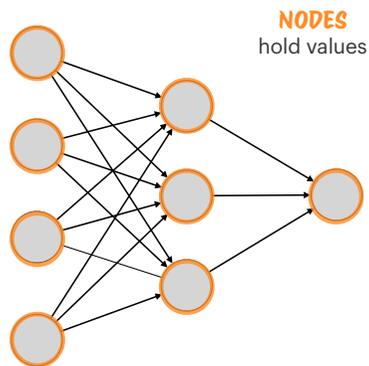
1

Terminology



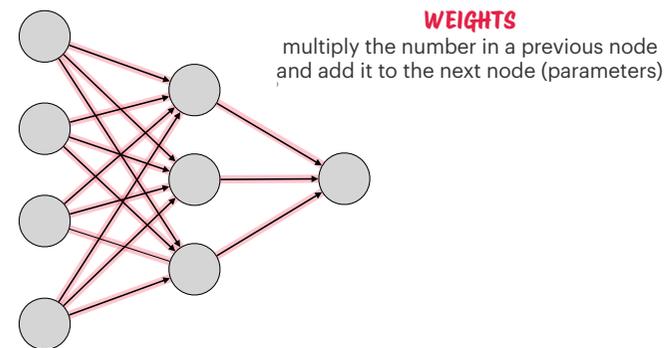
2

Terminology



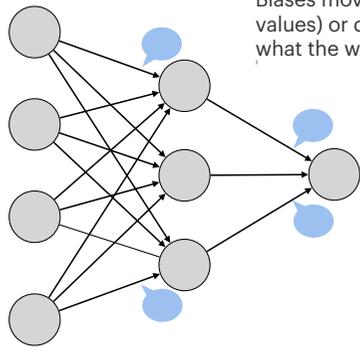
3

Terminology



4

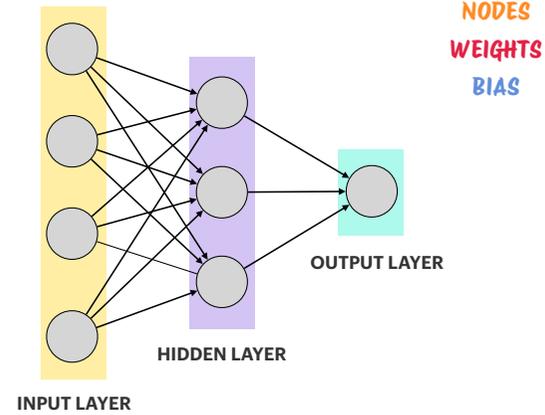
Terminology



BIAS
Biases move the value of a node up (for positive values) or down (for negative values) no matter what the weights and previous nodes' values were

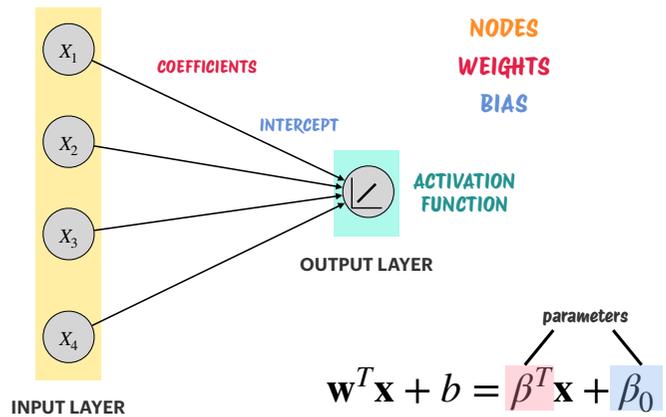
5

Terminology



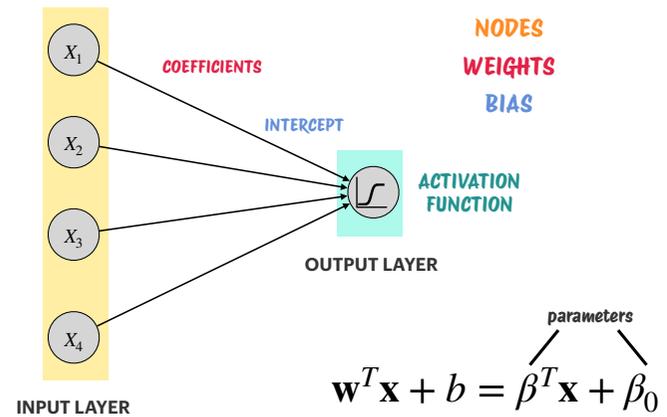
6

Linear Regression as a Neural Network



7

Logistic Regression as a Neural Network

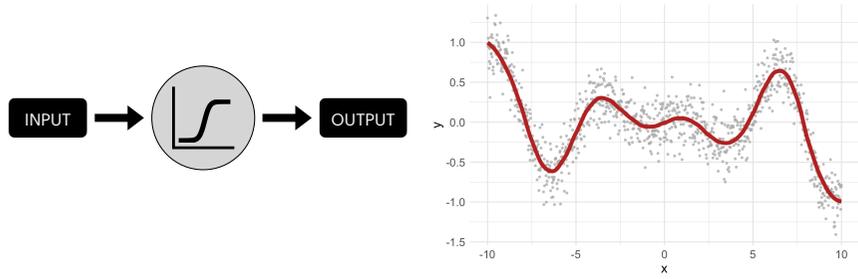


8

Activation Function

the purpose of an activation function is to add non-linearity to the neural network.

they determine whether a neuron should be activated based on its input



Without activation functions, the hidden layers would only apply linear transformations (weighted sums), and the entire network is collapsed into a single-layer model because the composition of two linear functions is a linear function itself

9

Activation Function

Activation Function	Equation	Range	Use Case
Sigmoid	$f(x) = \frac{1}{1 + e^{-x}}$	(0, 1)	Binary classification, hidden layers in small networks.
Tanh	$f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$	(-1, 1)	Hidden layers, encourages zero-centered outputs.
ReLU	$f(x) = \max(0, x)$	$[0, \infty)$	Default for most hidden layers, efficient computation.
Leaky ReLU	$f(x) = \max(0.01x, x)$	$(-\infty, \infty)$	Avoids dying ReLU problem, suitable for hidden layers.
Softmax	$f(x) = \frac{e^{x_j}}{\sum_j e^{x_j}}$	(0, 1)	Output layers in multi-class classification problems.
Swish	$f(x) = x \cdot \text{sigmoid}(x)$	$(-\infty, \infty)$	Recent innovation, smooth activation for better gradients.
GELU	$f(x) = 0.5x(1 + \tanh(\sqrt{2/\pi}(x + 0.044715x^3)))$	$(-\infty, \infty)$	High-performing activation in transformers.

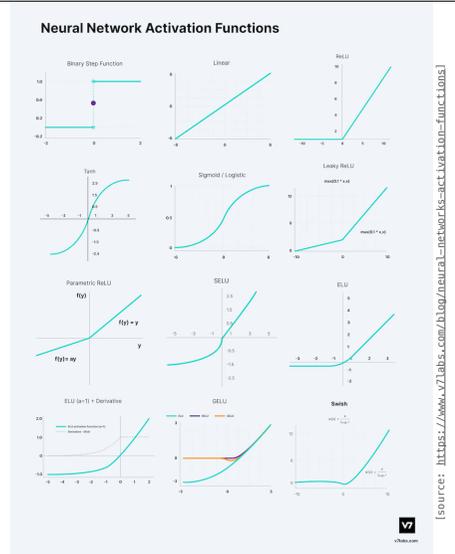
10

Activation Function

Examples:

- Regression — Linear Activation
- Binary Classification— Sigmoid/Logistic
- Multiclass Classification—Softmax
- Multilabel Classification—Sigmoid

⋮



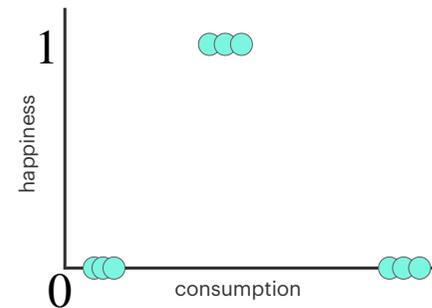
11

How Do Neural Networks Work?

a **VERY** simple example:

X = normalized chocolate consumption (0–1)

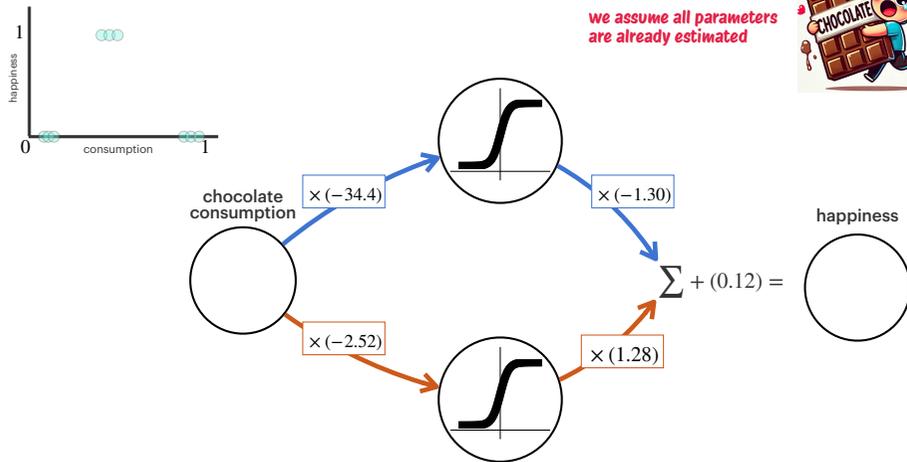
Y = happiness (0/1)



12

How Do Neural Networks Work?

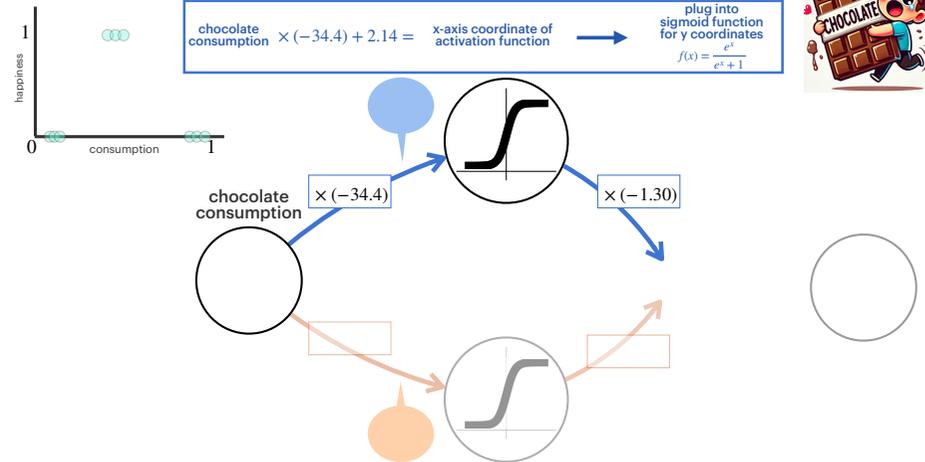
we assume all parameters are already estimated



13

How Do Neural Networks Work?

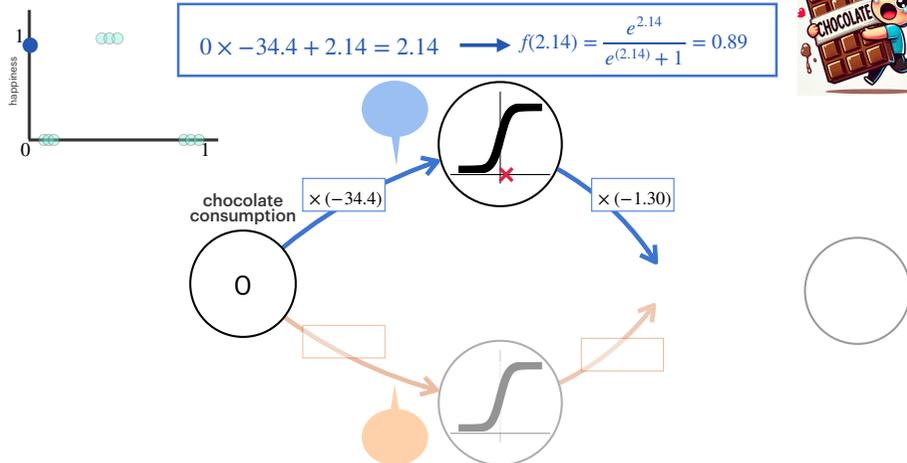
chocolate consumption $\times (-34.4) + 2.14 =$ x-axis coordinate of activation function \rightarrow plug into sigmoid function for y coordinates $f(x) = \frac{e^x}{e^x + 1}$



14

How Do Neural Networks Work?

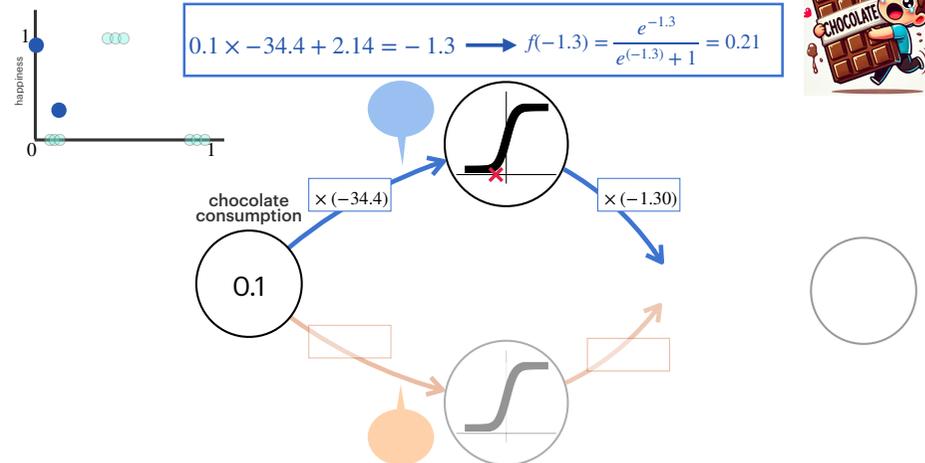
$0 \times -34.4 + 2.14 = 2.14 \rightarrow f(2.14) = \frac{e^{2.14}}{e^{2.14} + 1} = 0.89$



15

How Do Neural Networks Work?

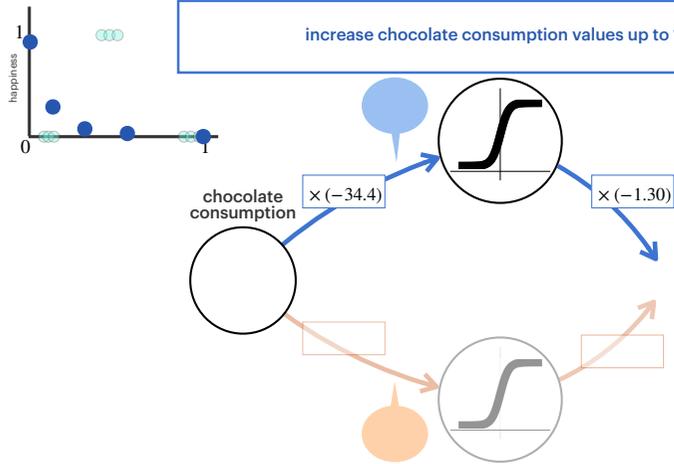
$0.1 \times -34.4 + 2.14 = -1.3 \rightarrow f(-1.3) = \frac{e^{-1.3}}{e^{-1.3} + 1} = 0.21$



16

How Do Neural Networks Work?

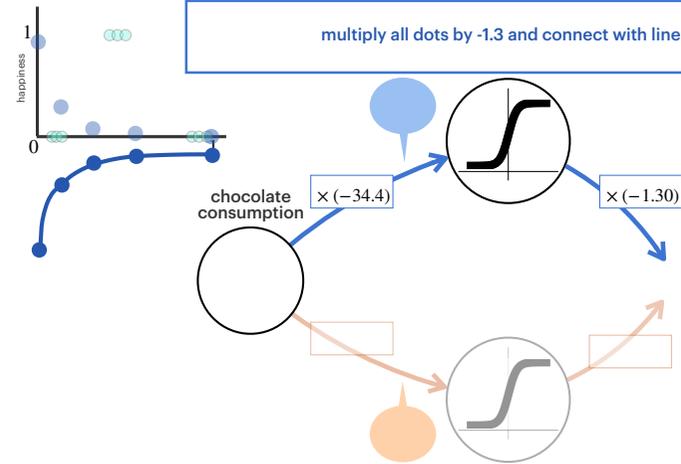
increase chocolate consumption values up to 1



17

How Do Neural Networks Work?

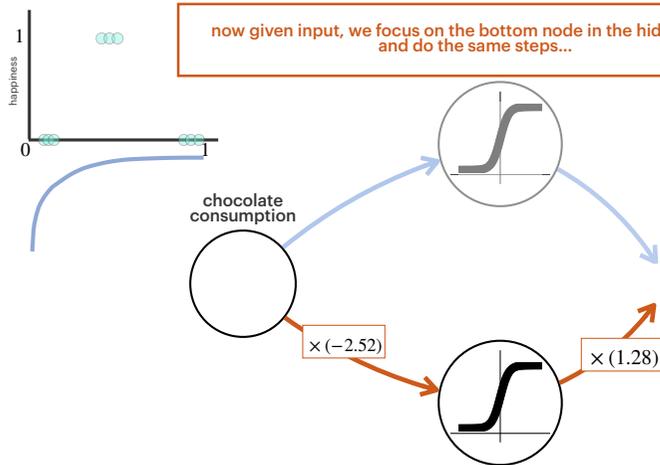
multiply all dots by -1.3 and connect with line



18

How Do Neural Networks Work?

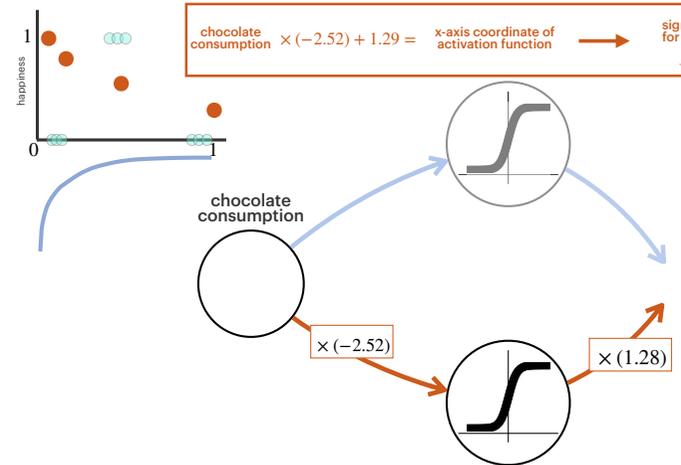
now given input, we focus on the bottom node in the hidden layer and do the same steps...



19

How Do Neural Networks Work?

chocolate consumption $\times (-2.52) + 1.29 =$ x-axis coordinate of activation function \rightarrow plug into sigmoid function for y coordinates
 $f(x) = \frac{e^x}{e^x + 1}$



20

How Do Neural Networks Work?

...multiply all dots by 1.28 and connect with line

chocolate consumption

$\times (-2.52)$

$\times (1.28)$

21

How Do Neural Networks Work?

next we sum the Y coordinates the blue and orange line...

chocolate consumption

$\times (-34.4)$

$\times (-2.52)$

$\times (-1.30)$

$\times (1.28)$

$\Sigma + (0.12) =$ happiness

22

How Do Neural Networks Work?

and then add 0.12 to get the final curve

chocolate consumption

$\times (-34.4)$

$\times (-2.52)$

$\times (-1.30)$

$\times (1.28)$

$\Sigma + (0.12) =$ happiness

23

Backpropagation

chocolate consumption

w_1

w_2

w_3

w_4

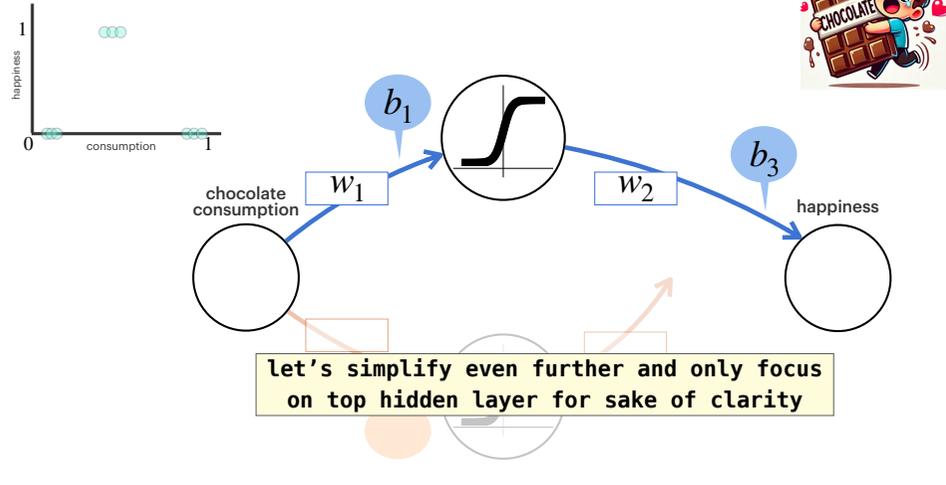
b_1

b_2

$\Sigma + b_3 =$ happiness

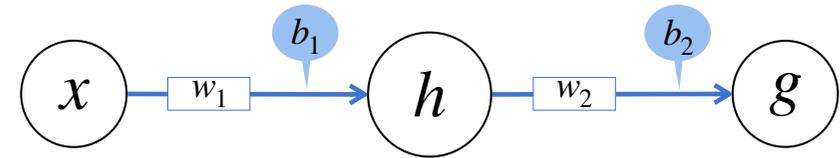
24

Backpropagation



25

Backpropagation



Chain rule:

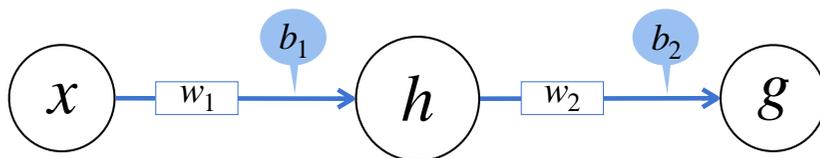
If we want to know how changing x affects $f(g(x))$ we first need to think about how changing x affects $g(x)$ and then how changing $g(x)$ affects $f(g(x))$:

$$\frac{\partial f(g(x))}{\partial x} = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial x}$$

the derivative of the outer function f evaluated at $g(x)$
multiplied by the derivative of the inner function $g(x)$

26

Backpropagation



$$h = w_1 \cdot x + b_1$$

$$g = w_2 \cdot x + b_2$$

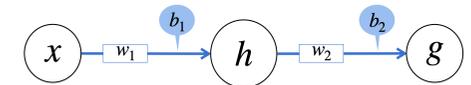
how close is g to our actual value?

Loss function (MSE):

$$\frac{1}{N} \sum_i (y_i - g_i)^2 \Rightarrow \frac{1}{N} \sum_i (y_i - \underbrace{(w_2 \cdot (w_1 \cdot x_i + b_1) + b_2)}_{= h})^2$$

27

Backpropagation



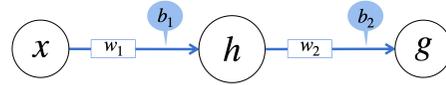
$$\frac{1}{N} \sum_i (y_i - g_i)^2 \Rightarrow \frac{1}{N} \sum_i (\underbrace{y_i}_{\text{actual}} - \underbrace{(w_2 \cdot (w_1 \cdot x_i + b_1) + b_2)}_{\text{predicted}})^2$$

this is what the gradient tells us!

- How change w_1 to reduce our loss?
- How change w_2 to reduce our loss?
- How change b_1 to reduce our loss?
- How change b_2 to reduce our loss?

28

Backpropagation



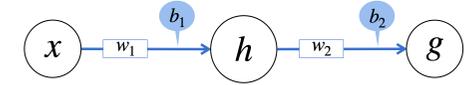
$$\frac{1}{N} \sum_i (y_i - g_i)^2 \implies \frac{1}{N} \sum_i \underbrace{(y_i}_{\text{actual}} - \underbrace{(w_2 \cdot (w_1 \cdot x_i + b_1) + b_2)}_{\text{predicted}})^2$$

$$\frac{\partial \text{Loss}}{\partial w_1} = \frac{\text{Loss}}{\partial g} \cdot \frac{\partial g}{\partial h} \cdot \frac{\partial h}{\partial w_1}$$

changing w_1 changes h , and changing h will change g ,
and changing g will change overall loss
 \implies we need the chain rule!

29

Backpropagation



$$\frac{1}{N} \sum_i (y_i - g_i)^2 \implies \frac{1}{N} \sum_i \underbrace{(y_i}_{\text{actual}} - \underbrace{(w_2 \cdot (w_1 \cdot x_i + b_1) + b_2)}_{\text{predicted}})^2$$

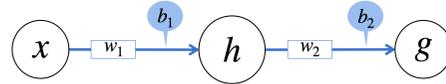
$$\frac{\partial \text{Loss}}{\partial w_1} = \underbrace{\frac{\text{Loss}}{\partial g}}_{-2(y_i - g_i)} \cdot \underbrace{\frac{\partial g}{\partial h}}_{w_2} \cdot \underbrace{\frac{\partial h}{\partial w_1}}_x$$

sum over all observations

this is the first part
of our gradient

30

Backpropagation



$$\frac{1}{N} \sum_i (y_i - g_i)^2 \implies \frac{1}{N} \sum_i \underbrace{(y_i}_{\text{actual}} - \underbrace{(w_2 \cdot (w_1 \cdot x_i + b_1) + b_2)}_{\text{predicted}})^2$$

next part of gradient
is with respect to our
second parameter

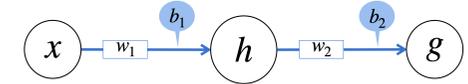
$$\frac{\partial \text{Loss}}{\partial b_1} = \frac{\text{Loss}}{\partial g} \cdot \frac{\partial g}{\partial h} \cdot \frac{\partial h}{\partial b_1}$$

$$\underbrace{-2(y_i - g_i) w_2}_{\text{sum over all observations}} \cdot 1$$

and so forth for all parameters...
this is backpropagation

31

Backpropagation



$$\frac{1}{N} \sum_i (y_i - g_i)^2 \implies \frac{1}{N} \sum_i \underbrace{(y_i}_{\text{actual}} - \underbrace{(w_2 \cdot (w_1 \cdot x_i + b_1) + b_2)}_{\text{predicted}})^2$$

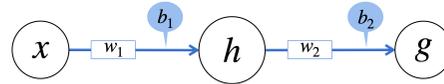
$$\frac{\partial L}{\partial w_1} = -\frac{2}{N} \sum_{i=1}^N (y_i - g_i) w_2 x_i \quad \frac{\partial L}{\partial w_2} = -\frac{2}{N} \sum_{i=1}^N (y_i - g_i) h_i$$

$$\frac{\partial L}{\partial b_1} = -\frac{2}{N} \sum_{i=1}^N (y_i - g_i) w_2 \quad \frac{\partial L}{\partial b_2} = -\frac{2}{N} \sum_{i=1}^N (y_i - g_i)$$

and so forth for all parameters...
this is backpropagation

32

Backpropagation



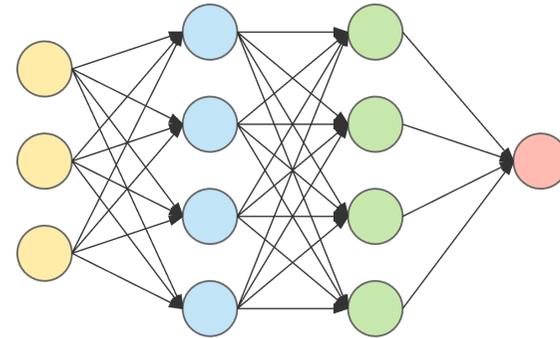
$$\frac{1}{N} \sum_i (y_i - g_i)^2 \implies \frac{1}{N} \sum_i \underbrace{(y_i}_{\text{actual}} - \underbrace{(w_2 \cdot (w_1 \cdot x_i + b_1) + b_2)}_{\text{predicted}})^2$$

$$\nabla_{\theta} L = \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial b_1} \\ \frac{\partial L}{\partial w_2} \\ \frac{\partial L}{\partial b_2} \end{bmatrix} = -\frac{2}{N} \sum_{i=1}^N \begin{bmatrix} e_i w_2 x_i \\ e_i w_2 \\ e_i h_i \\ e_i \end{bmatrix} \quad \text{where } e_i = y_i - g_i$$

and so forth for all parameters...
this is backpropagation

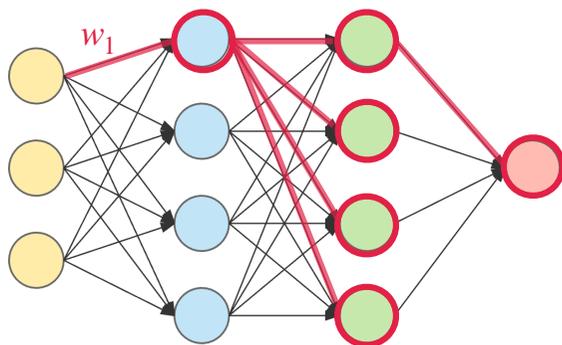
33

Butterfly Effect



34

Butterfly Effect



35

Training Process

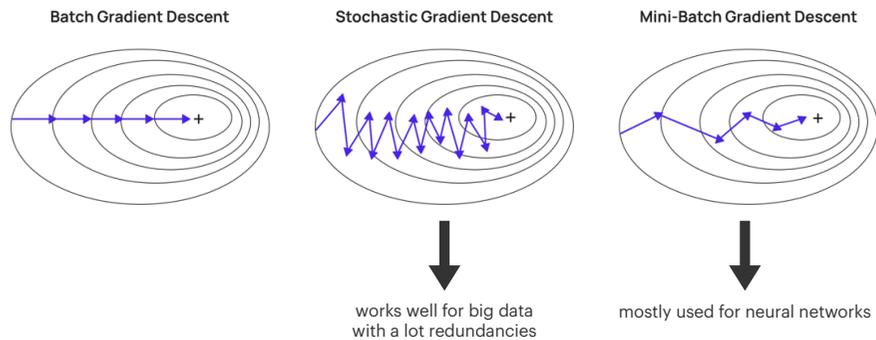
	Forward Pass	Backward Pass
Direction	Input → Hidden Layers → Output	Output → Hidden Layers → Input
Purpose	Calculate output and loss	Update weights and biases to reduce loss
Mathematics	Linear transformations and activations	Gradients and chain rule
Focus	Prediction	Learning (optimization)

what are **Epochs**?

- Refers to one complete pass through the entire training dataset by the model. During this
 1. the model sees every single example in the training dataset once.
 2. It tries to learn by adjusting its weights through forward passes (making predictions) and backward passes (updating weights using the errors).
- More epochs = more practice = better learning (up to a point; too many epochs can lead to overfitting)

36

Versions of Gradient Descent



37

Problems with Gradient Descent

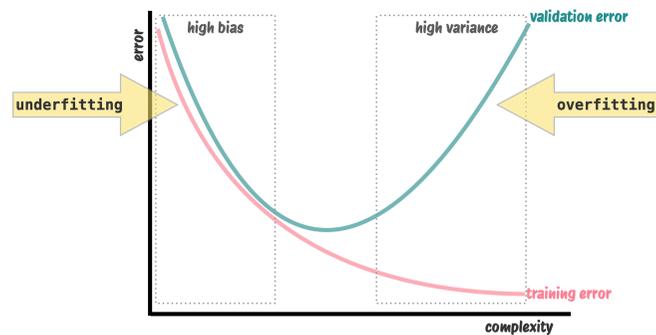
- Choosing a learning rate is hard
- Setting learning rate schedules ahead of time is hard
- Using the same learning rate for all parameters
- Local minima and saddle points
- Variations exist to overcome some of these problems, e.g.
 - **Momentum:**
 - Momentum allows us to “build up speed” as we go downhill
 - It uses a moving average
 - Momentum often helps us converge faster
 - Momentum allows us to escape (some) local minima

38

Regularization

with great complexity comes great regularization

Bias Variance Trade-Off



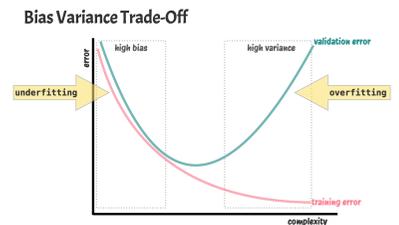
39

Regularization

with great complexity comes great regularization

Penalization

- best for smaller networks
- added to each layer
- L1, L2, L1 + L2



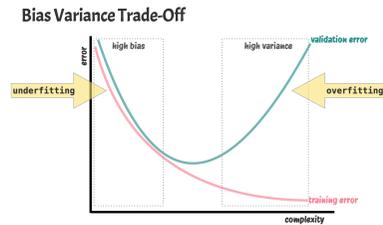
40

Regularization

with great complexity comes great regularization

Bagging

- we can use bagging to get different subsamples of there data to train different neural nets on
- average results together
- computationally expensive



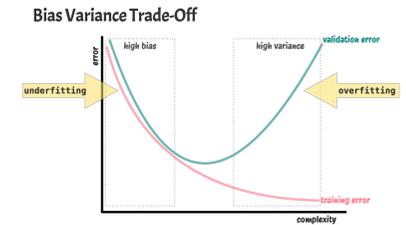
41

Regularization

with great complexity comes great regularization

Dropout

- randomly selects nodes in each of the layers and removes them/zeroes them out
- applied to each layer
- networks cannot over-rely on some connections



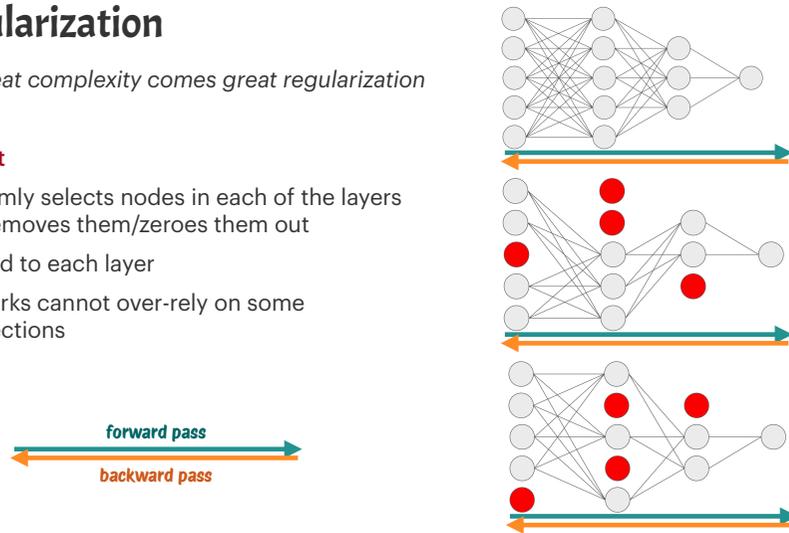
42

Regularization

with great complexity comes great regularization

Dropout

- randomly selects nodes in each of the layers and removes them/zeroes them out
- applied to each layer
- networks cannot over-rely on some connections



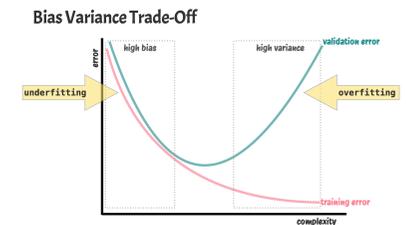
43

Regularization

with great complexity comes great regularization

Early Stopping

- while training, when our validation error stops improving or gets worse, then stop the training process
- "patience" for stopping = number of iterations where our validation error can remain the same or get worse



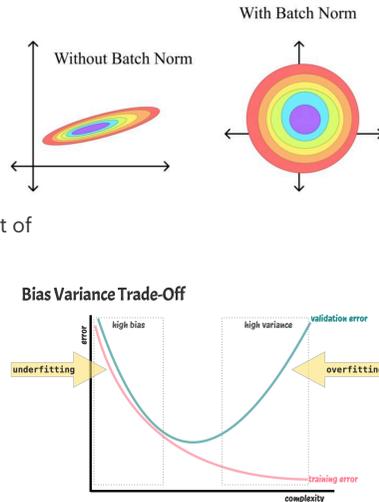
44

Regularization

with great complexity comes great regularization

Batch Normalization

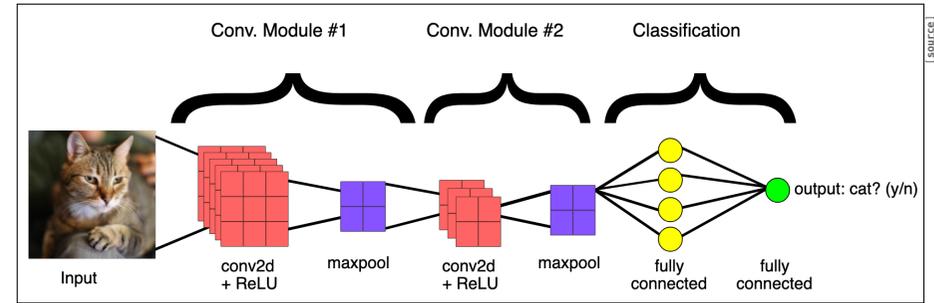
- Stabilizes learning
rescales each layer's activations to have a consistent mean and variance during training and lessens impact of initial weights
- Speeds up and regularizes training
allows faster convergence and often reduces overfitting



45

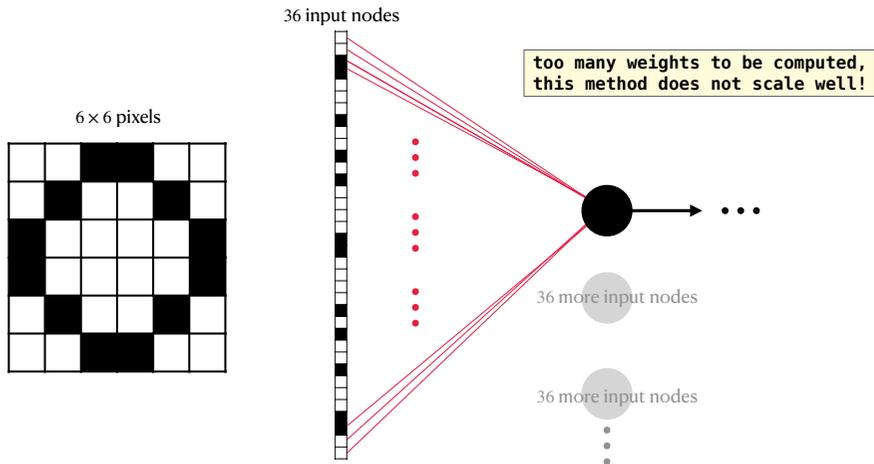
Image Classification with Convolutional Networks

neural networks with weight sharing



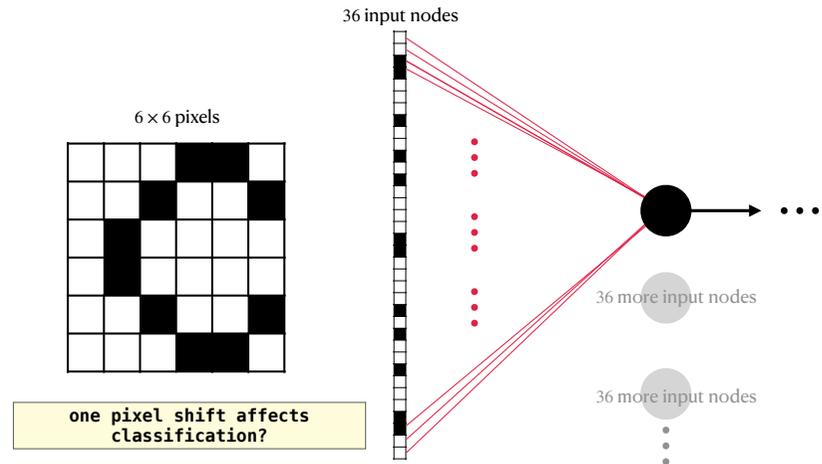
46

Image Classification with Convolutional Networks



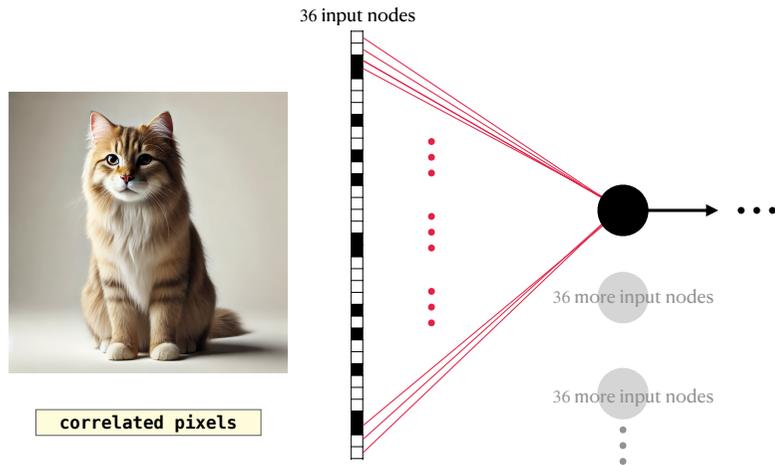
47

Image Classification with Convolutional Networks



48

Image Classification with Convolutional Networks



49

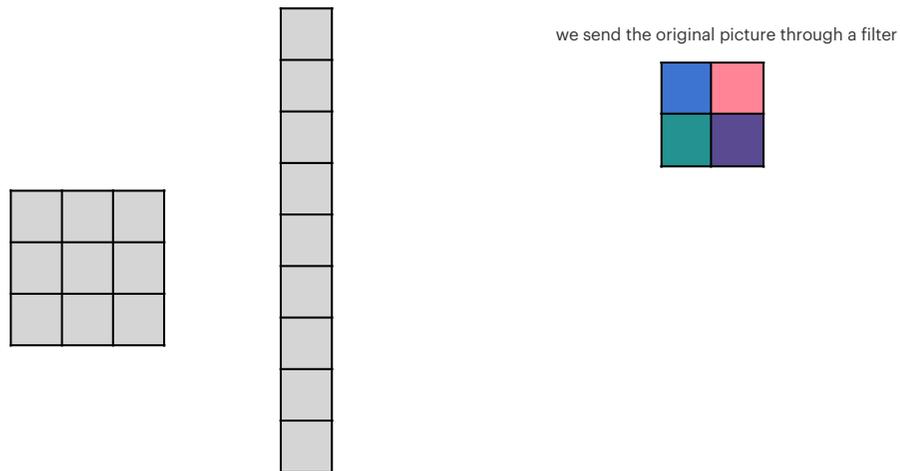
Image Classification with Convolutional Networks

- Convolutional networks will
- reduces number of input nodes
 - tolerate pixel shifts in image
 - take advantage of pixel correlations



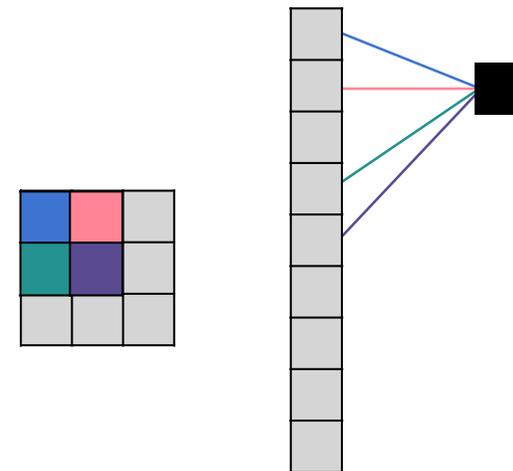
50

Image Classification with Convolutional Networks



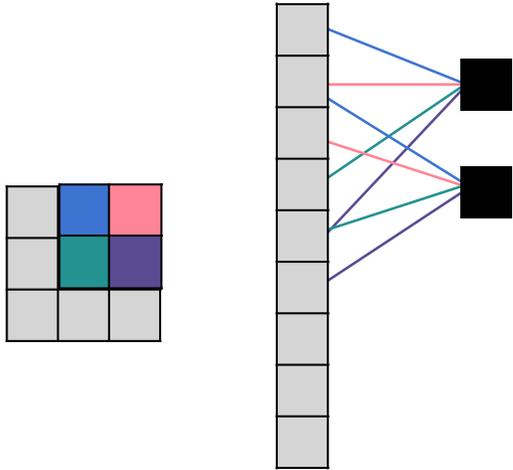
51

Image Classification with Convolutional Networks



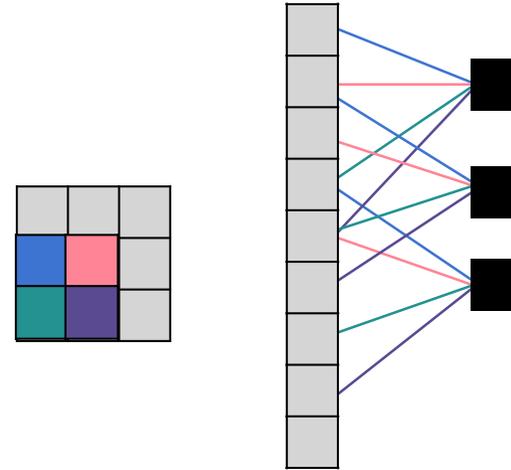
52

Image Classification with Convolutional Networks



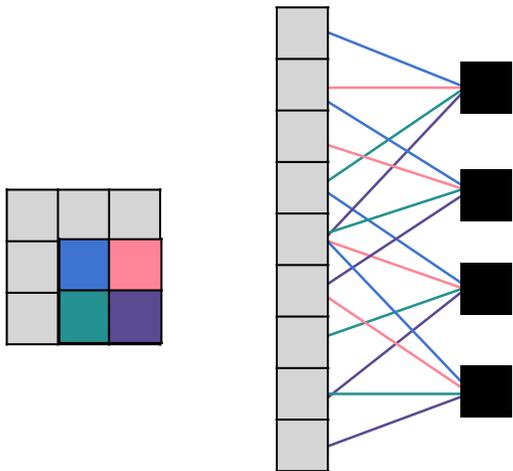
53

Image Classification with Convolutional Networks



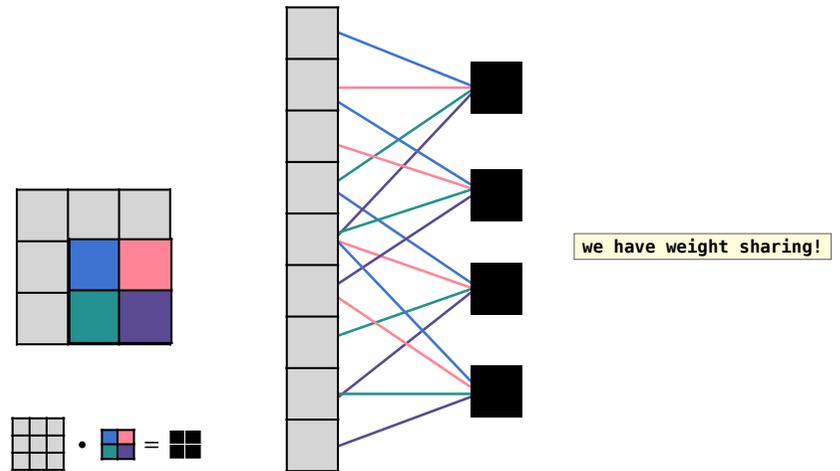
54

Image Classification with Convolutional Networks



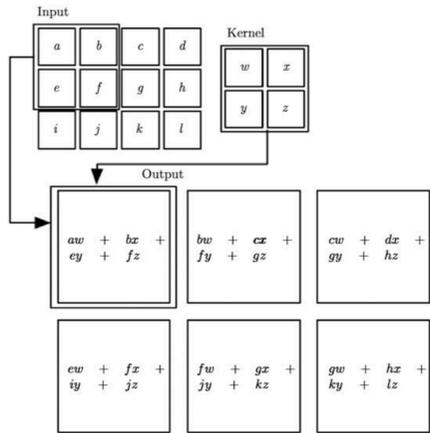
55

Image Classification with Convolutional Networks



56

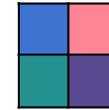
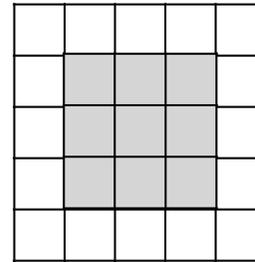
Image Classification with Convolutional Networks



From Deep Learning by Goodfellow, Bengio and Courville

57

Image Classification with Convolutional Networks

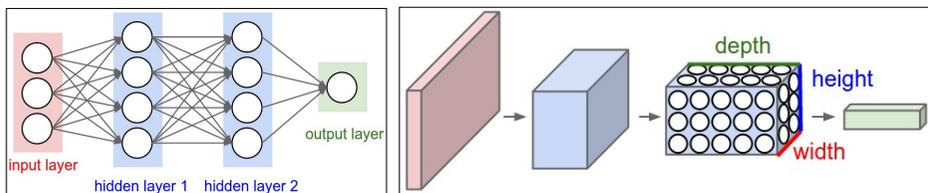


note: the resulting image will be $(n-1) \times (n-1)$ to avoid that and keep original dimensions we can use *padding*

58

Image Classification with Convolutional Networks

neural networks with weight sharing and **3D activations**



59

Image Classification with Convolutional Networks

strides

strides determine how much the convolution operation simplifies or compresses the input data

- examples

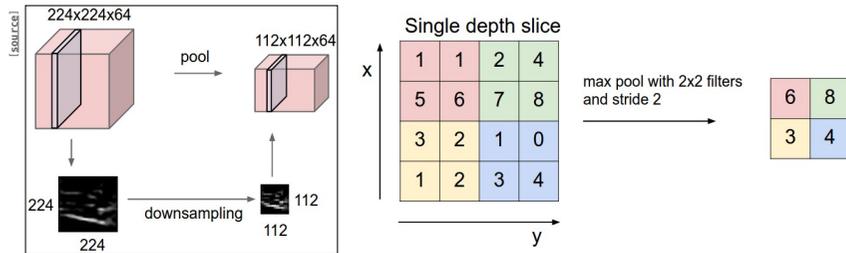
- Stride = 1 (default): the kernel moves one pixel at a time horizontally or vertically. High degree of overlap receptive fields but usually retains more spatial details.
- Stride > 1: The kernel skips pixels as it moves, reducing overlap between receptive fields. This leads to a smaller output size and down-sampling of the input image.

60

Image Classification with Convolutional Networks

pooling

- Downsample feature maps
- Max pooling works best because when looking for a feature it is better to look at the maximal presence rather than the average presence
- It's best practice to do un-strided convolutions then downsample with maxpooling rather than using strides to downsample

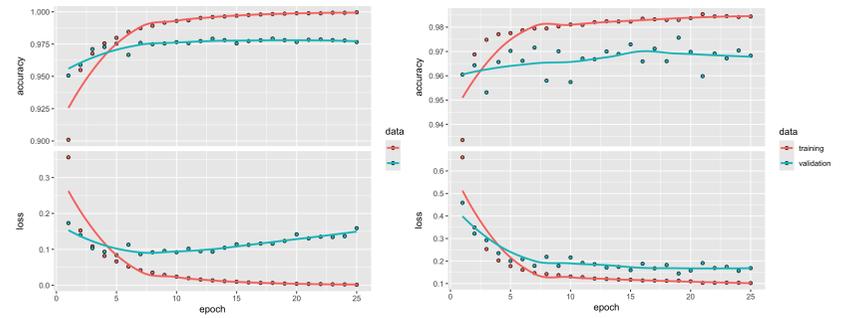


61

This Week's Practical Neural Networks

```

0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9 9 9
    
```



62