

# Convolutional Neural Networks (CNN)

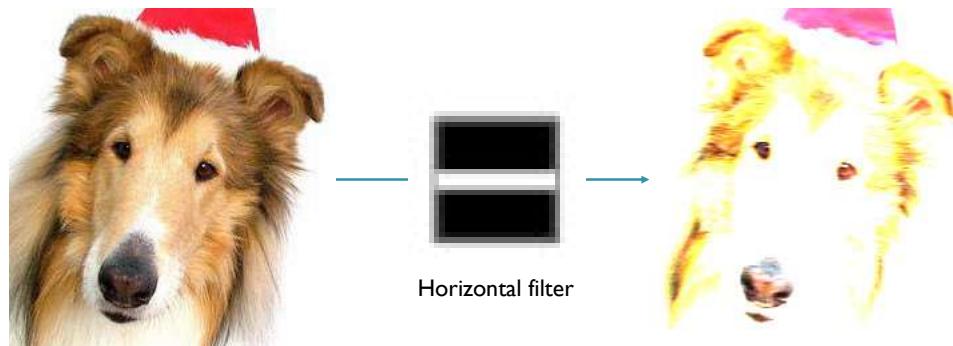
©2018 Peter V. Henstock and David Dowey

A picture says a 1000 words



©2018 Peter V. Henstock and David Dowey

## Feature selection



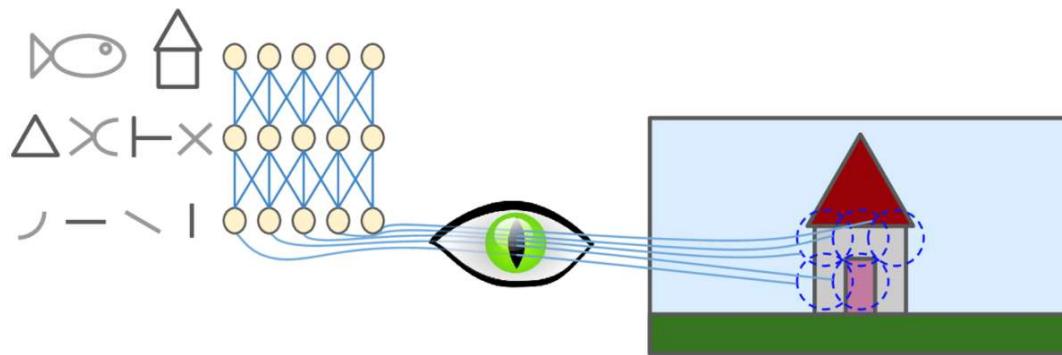
Horizontal filter

©2018 Peter V. Henstock and David Dowey

## Why Convolutional Neural Networks?

Neural network analogy with human vision system and the architecture of the visual cortex

- Local receptive field: limited regions (e.g. 5 neurons/layer)
- Complex patterns are combinations of lower-level patterns



from the book Hands on Machine Learning with Scikit-Learn and TensorFlow, © Géron, 2017

©2018 Peter V. Henstock and David Dowey

## Why Convolutional Neural Networks?

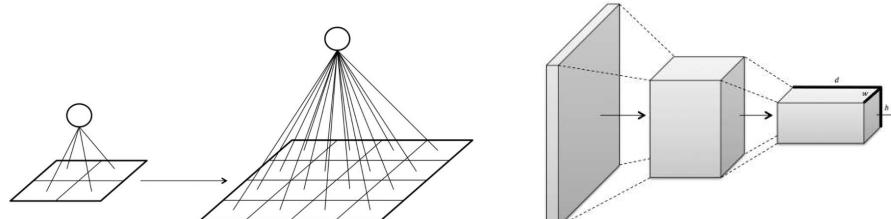
Ideas have fed into & driven NN research

- One idea is effective scaling
  - Not all neurons fire from an image
  - Distribution based on features
- Second idea is NN can recognize faces in different lighting/angles/obscuration
- Third idea is that it is the combination of features that define recognition and works even for partially obscured images
- Fourth idea is that there are hierarchical relationships that can help learning

©2018 Peter V. Henstock and David Dowey

## Why Convolutional Neural Networks?

- MNIST uses 28x28 images → 784 weights
- Scaling up to 256x256 would kill them
- Need a scalable solution
- 2012 ImageNet benchmark → 25.7% error
  - Alex Krizhevsky (Geoffrey Hinton's Lab)
  - AlexNet used 1<sup>st</sup> CNN; later → 16% error rate



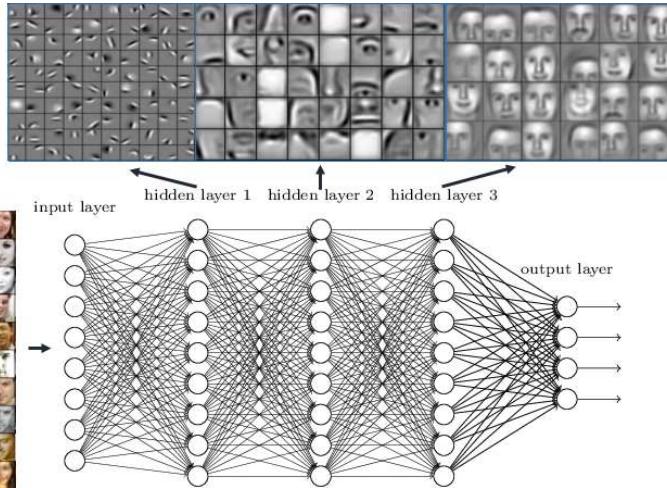
from the book Fundamentals of Deep Learning, © Buduma, 2017

©2018 Peter V. Henstock and David Dowey

## What Deep Learning is Doing

- <http://www.rsipvision.com/exploring-deep-learning/>

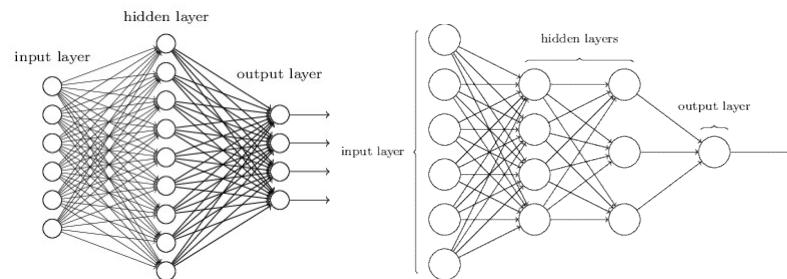
Deep neural networks learn hierarchical feature representations



©2018 Peter V. Henstock and David Dowey

## Layers vs. Hidden Layer Width

- Both might have the same number of nodes but which to choose?



- Generally #connections is the killer
- Wide hidden layer has more connections

©2018 Peter V. Henstock and David Dowey

## Observation

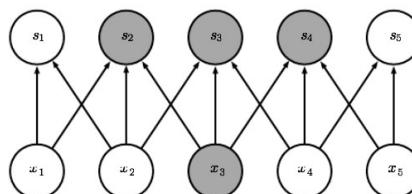
- Fully-connected layer: Connection goes from every node of previous layer to every node of next layer
- What if we relaxed that assumption?
- Brain uses regions of connected neurons
- NN with small regions of connectivity
  - Convolutional Neural Network
  - Also called “local networks”
  - Also called “locally connected networks”

©2018 Peter V. Henstock and David Dowey

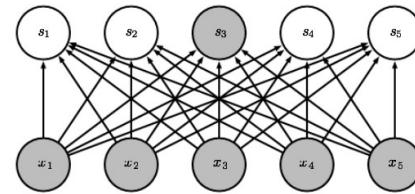
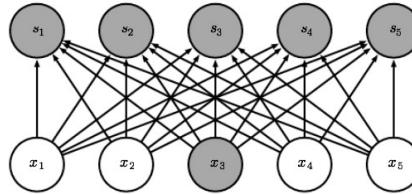
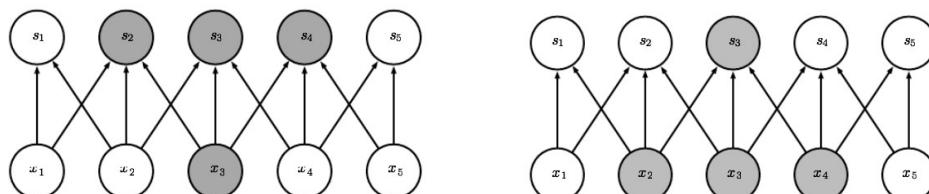
## Densely connected vs. sparse connectivity

- kernel smaller than input → detect small, meaningful features
- store fewer weights, use less memory, improve statistical efficiency

view from below



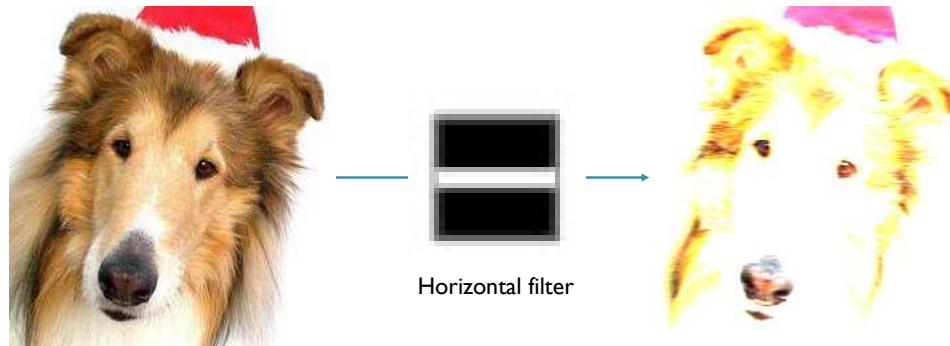
view from above



from the book Deep Learning, © Goodfellow et al., 2016

©2018 Peter V. Henstock and David Dowey

## Weight Sharing

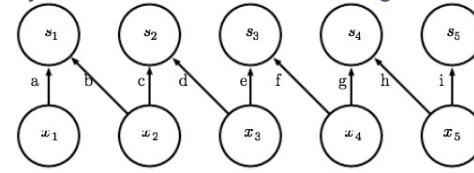


©2018 Peter V. Henstock and David Dowey

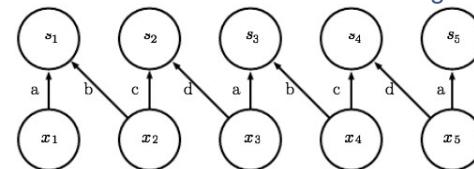
## Weight Sharing

- a locally connected (sparse) network is not obliged to share weights
- in tiled convolutions, kernels rotate – different filters are applied to neighboring neurons
- sharing weights in a traditional Convolutional Neural Network has only one kernel that passes everywhere

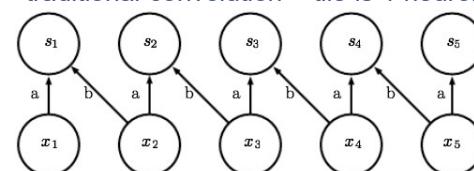
locally connected network – no weight sharing



tiled – 2 neurons have different weights



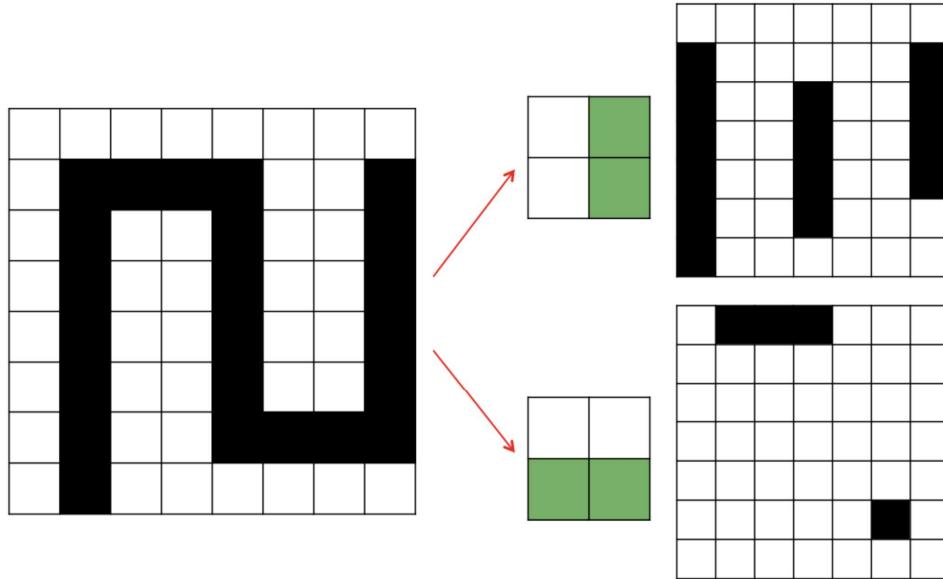
traditional convolution – tile is 1 neuron



from the book Deep Learning, © Goodfellow et al., 2016

©2018 Peter V. Henstock and David Dowey

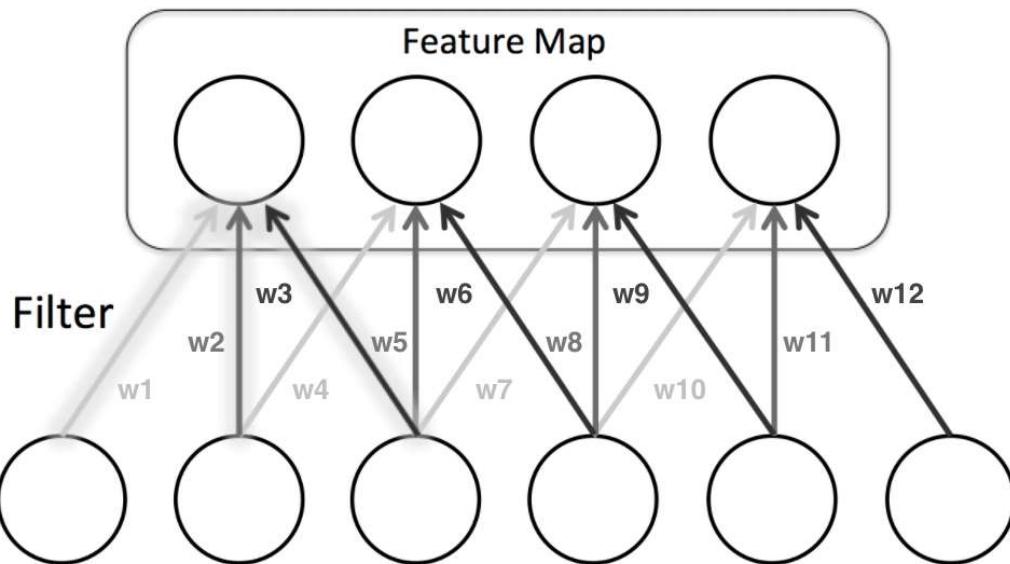
## Weight Sharing



from the book Fundamentals of Deep Learning, © Buduma, 2017

©2018 Peter V. Henstock and David Dowey

## Weight Sharing



from the book Fundamentals of Deep Learning, © Buduma, 2017

©2018 Peter V. Henstock and David Dowey

## Weight sharing

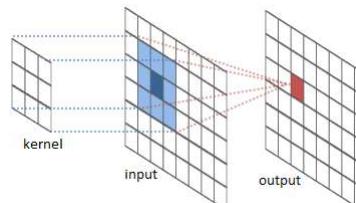
- Learning weights but shared (i.e. equal)
- Shifting by 3 here each time:
  - $w_1 = w_4 = w_7 = w_{10}$
  - $w_2 = w_5 = w_8 = w_{11}$
  - $w_3 = w_6 = w_9 = w_{12}$

©2018 Peter V. Henstock and David Dowey

## Convolution in 2D

- Use 2D mask
- Apply to all locations
  - Multiply mask
  - Add products
  - Store result to center pixel
- Goal: Find a head

<https://medium.com/@ageitgey/machine-learning-is-fun-part-3-deep-learning-and-convolutional-neural-networks-f40359318721#nltv14pk4>



<http://stackoverflow.com/questions/15356153/how-do-convolution-matrices-work>



©2018 Peter V. Henstock and David Dowey

## Convolution

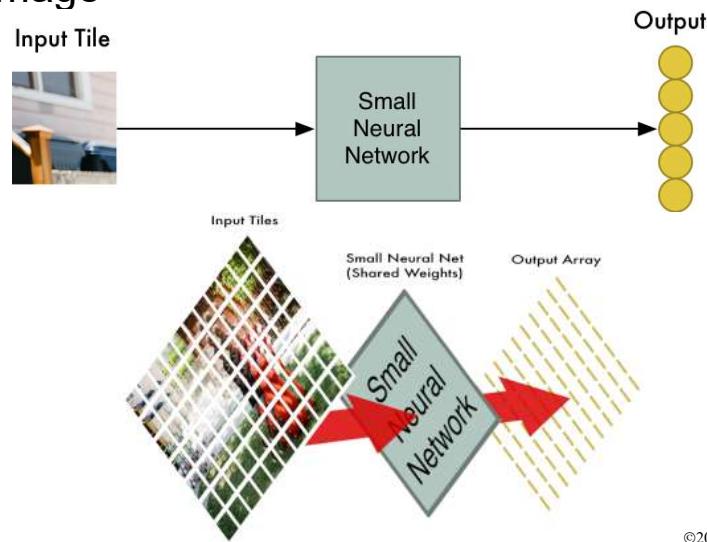


- We could have shifted every pixel
- Typically take a “stride” that will overlap
  - Stride = distance of the shift

©2018 Peter V. Henstock and David Dowey

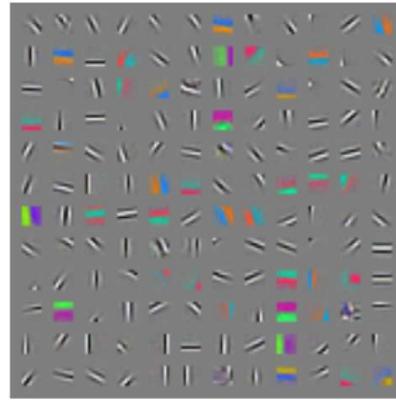
## Feature Reduction

- Train neural network to produce outputs for each subimage



©2018 Peter V. Henstock and David Dowey

## Weight sharing



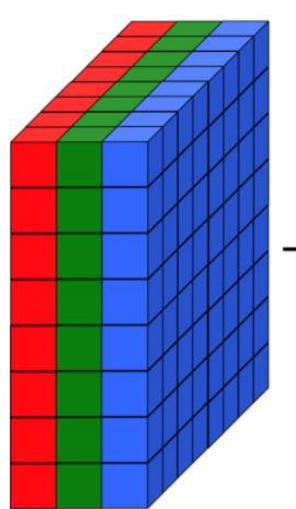
Convolutional kernels (filters) learned in a CNN

Elements in one of the feature maps are combinations of filter multiplications over all slices of the previous layer's feature maps

from the book Deep Learning, © Goodfellow et al., 2016

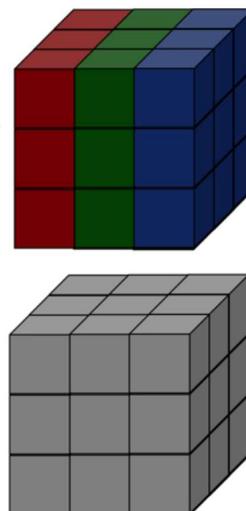
©2018 Peter V. Henstock and David Dowey

## Weight Sharing



RGB Image

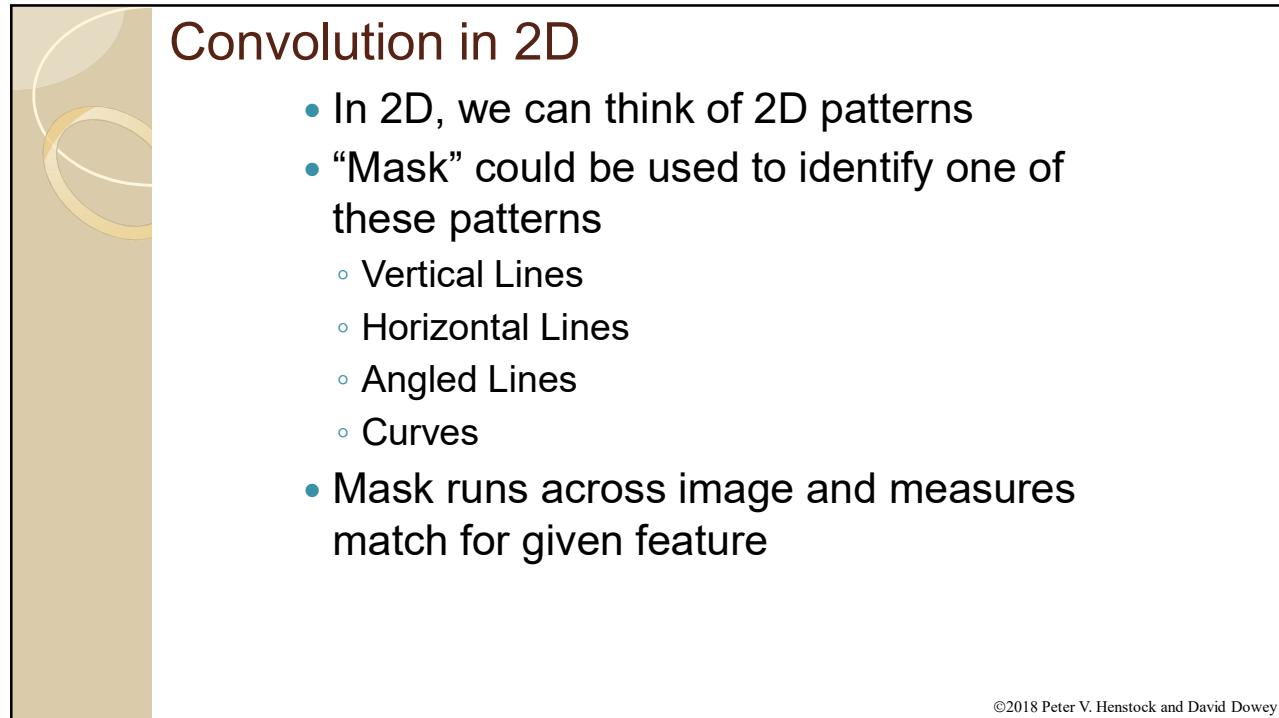
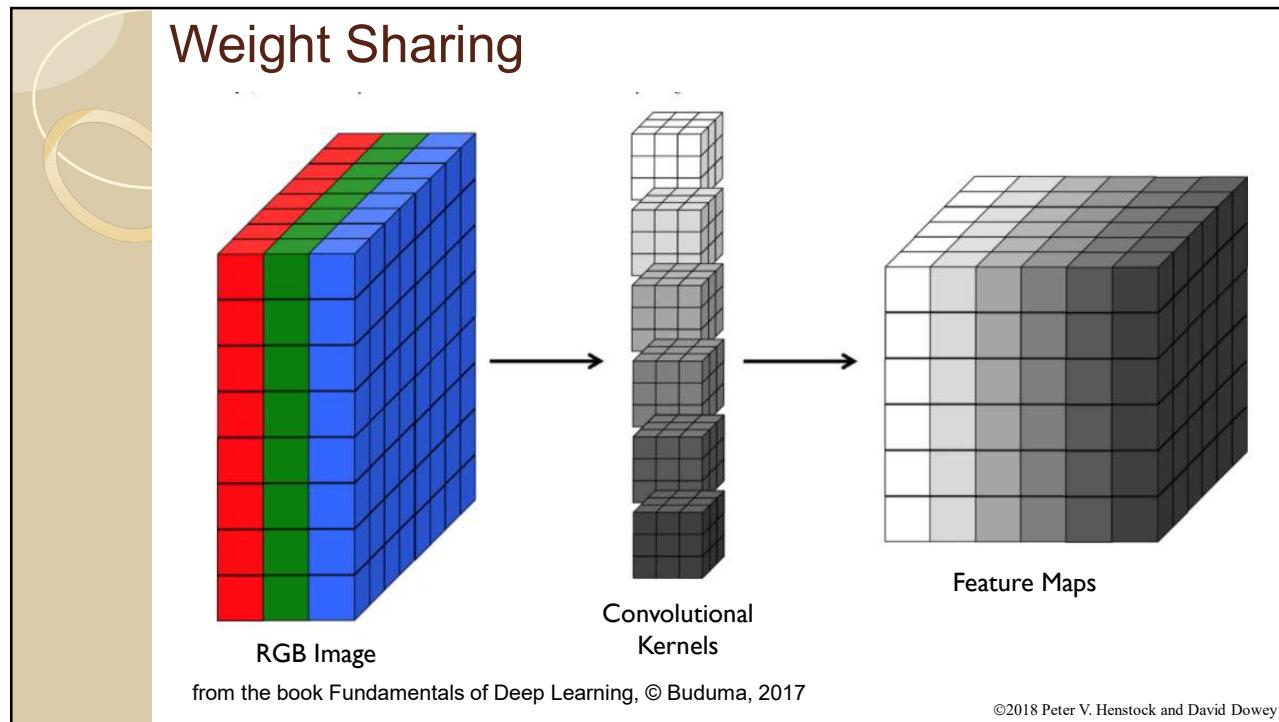
from the book Fundamentals of Deep Learning, © Buduma, 2017



Convolutional  
Kernels

?

©2018 Peter V. Henstock and David Dowey

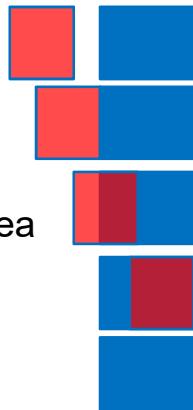
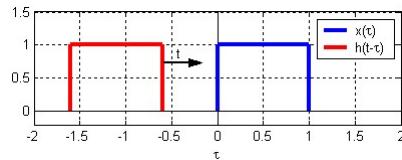


## The mechanics of convolution

- Signal processing mechanism for matching filter against a stream

$$[f * g](t) \equiv \int_{-\infty}^{\infty} f(\tau) g(t - \tau) d\tau,$$

- Video <http://mathworld.wolfram.com/Convolution.html>



- Convolution gives overlap area function of shift
- What shape will it be here?

©2018 Peter V. Henstock and David Dowey

## 2D convolution arithmetic

- A guide to convolution arithmetic in deep learning
- <https://arxiv.org/pdf/1603.07285.pdf>

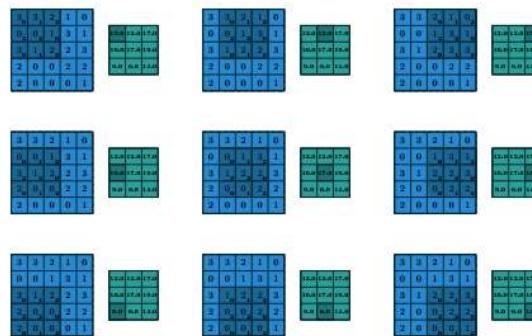
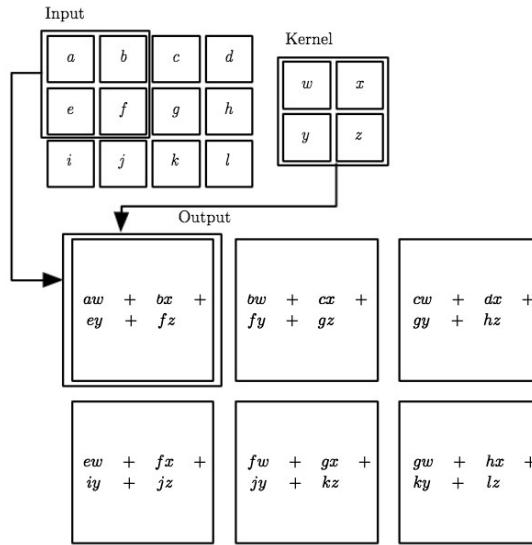


Figure 1.1: Computing the output values of a discrete convolution.

©2018 Peter V. Henstock and David Dowey

## 2D convolution arithmetic



from the book Deep Learning, © Goodfellow et al., 2016

- Input is 3 x 4 image
- kernel is 2 x 2 (patch)
- Output is called a "feature map": it is 2 x 3 and composed of the sum of element by element multiplications
- "VALID" padding: apply the kernel only entirely within the image – no padding

©2018 Peter V. Henstock and David Dowey

## 2D convolution arithmetic

Input 5 x 5

1	2	1	1	1
1	0	0	0	0
1	0	0	0	0
5	2	5	4	5
3	5	5	5	3

Kernel 2 x 2

-1	-1
1	1

Top left:  $-1*1 + -1*2 + 1*1 + 1*0 = -1$

Feature Map 4 x 4

-1	-3	-2	-2
0	0	0	0
6	7	9	9
1	3	1	-1

$$s(t) = (x * w)(t) = \sum_{a=-\infty}^{\infty} x(a)w(t-a).$$

from the book Deep Learning, © Goodfellow et al., 2016

©2018 Peter V. Henstock and David Dowey

## 2D convolution arithmetic

Online demos for convolution

- Setosa:

<http://setosa.io/ev/image-kernels/>

- Stanford CS231n:

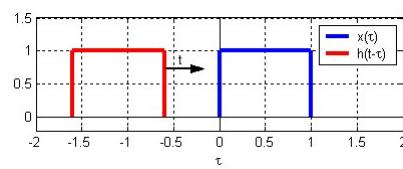
<https://cs231n.github.io/assets/conv-demo/index.html>

And a nice blog from Chris Olah:

<http://colah.github.io/posts/2014-07-Understanding-Convolutions/>

©2018 Peter V. Henstock and David Dowey

## What about the weights?



-1	0	+1
-2	0	+2
-1	0	+1

Gx

+1	+2	+1
0	0	0
-1	-2	-1

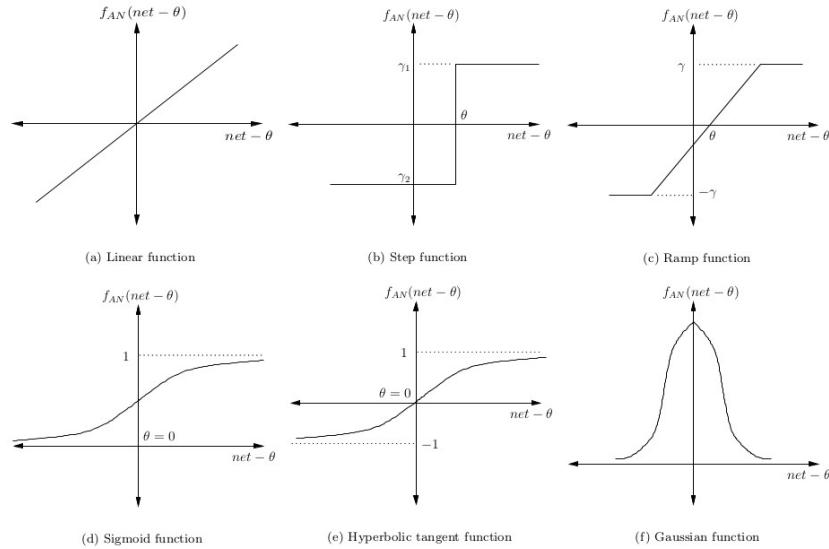
Gy

- Have masks running across the image focusing on a local region
- Where are the weights in all this?
- Is this a linear or non-linear operation?

©2018 Peter V. Henstock and David Dowey

## Activation Functions

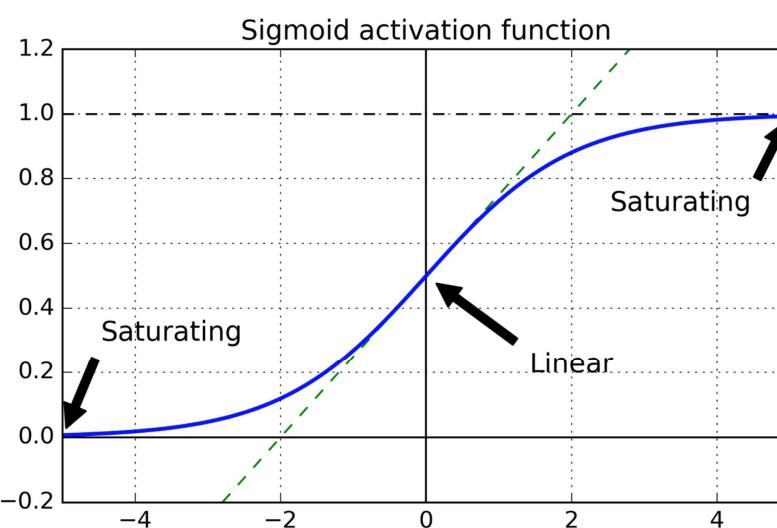
- <http://www.turingfinance.com/misconceptions-about-neural-networks/>



©2018 Peter V. Henstock and David Dowey

## Activation Functions

- Could Mother Nature have got it wrong?

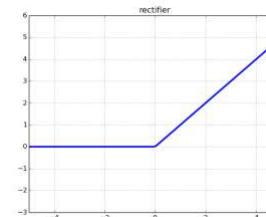


from the book Hands on Machine Learning with Scikit-Learn and TensorFlow, © Géron, 2017

©2018 Peter V. Henstock and David Dowey

## Sigmoid Calculation → ReLU

- Most of the past few decades have focused on sigmoid calculation
- Then we turned to ReLU and its variations
  - “Rectified Linear Unit”
- Softplus
  - Smooth version of this
  - $f(x) = \ln(1+\exp(x))$
- Why ReLU?



©2018 Peter V. Henstock and David Dowey

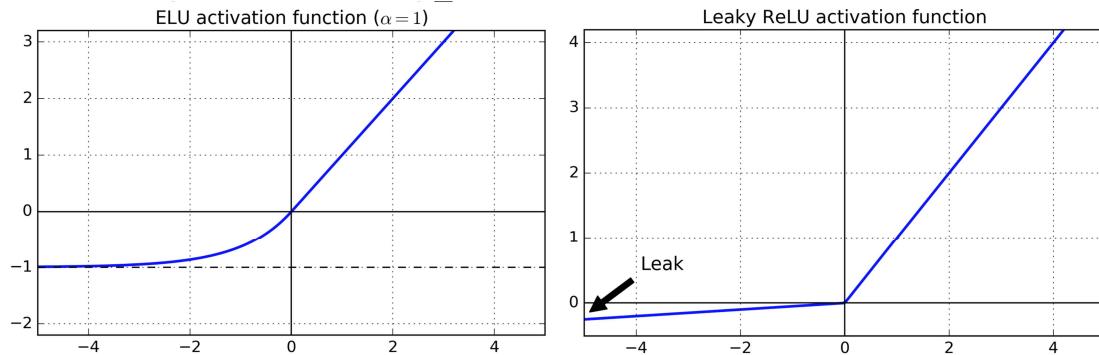
## Why ReLU

- Faster to calculate – not a concern for the brain
- Sigmoid gradients at plateaus
  - Tend to be small numbers not 0
  - Small numbers require calculation; not 0
    - With zeros, get sparser data and “it works”
  - ReLU accelerates stochastic gradient descent possibly since it doesn’t saturate
- Some of the nodes die and never activate – always 0 vanishing/exploding gradients problem

©2018 Peter V. Henstock and David Dowey

## Nonsaturating Activation Functions

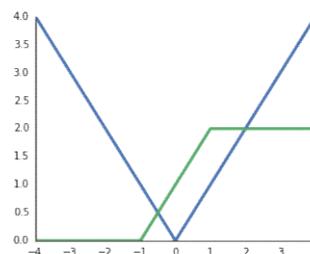
Leaky ReLU, Randomized Leaky ReLU (RReLU), Parametric Leaky (PReLU), Exponential Linear Unit (ELU), SWISH  
A long coma is sometimes better than death



from the book Hands on Machine Learning with Scikit-Learn and TensorFlow, © Géron, 2017  
©2018 Peter V. Henstock and David Dowey

## ReLU

- One reason for multiple layers was to include non-linearities
- Does the ReLU provide non-linearity?
- Answer:
  - 1) can across multiple layers
  - 2) doesn't have to be the only activation

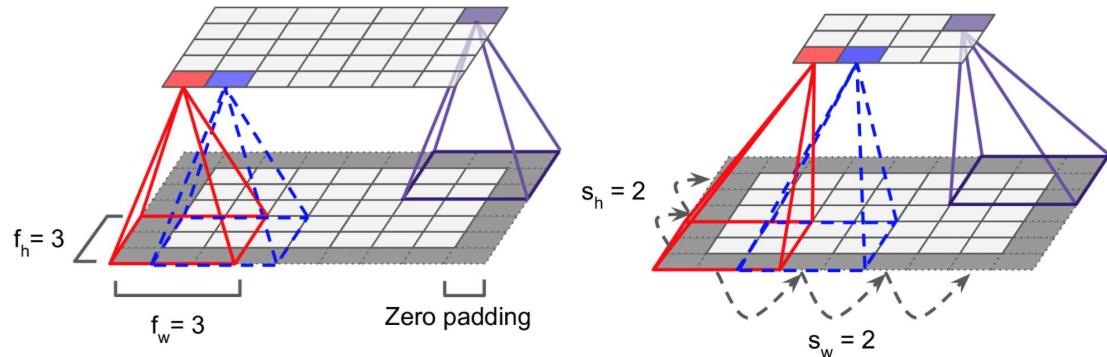


<http://stats.stackexchange.com/questions/141960/deep-neural-nets-relus-removing-non-linearity> ©2018 Peter V. Henstock and David Dowey

## Padding and stride

$3 \times 3$  patch with "SAME" (i.e. with zeros) padding on top and bottom

A stride of 2 – in both the height and width-wise directions of the image



from the book Hands on Machine Learning with Scikit-Learn and TensorFlow, © Géron, 2017

©2018 Peter V. Henstock and David Dowey

## Padding and stride

**padding="VALID"**  
(i.e., without padding)

Ignored

**padding="SAME"**  
(i.e., with zero padding)

from the book Hands on Machine Learning with Scikit-Learn and TensorFlow, © Géron, 2017

©2018 Peter V. Henstock and David Dowey

Figure 13-7. Padding options—input width: 13, filter width: 6, stride: 5

## Stacking multiple feature maps

- Start with 3 channels (RGB);
- Each convolutional layer  $l$  has  $n$  filters (kernels), each of them learning a characteristic and producing its own feature map  $k \rightarrow k = 1 \dots n$  feature maps

*Equation 13-1. Computing the output of a neuron in a convolutional layer*

$$z_{i,j,k} = b_k + \sum_{u=0}^{f_h-1} \sum_{v=0}^{f_w-1} \sum_{k'=0}^{f_n-1} x_{i',j',k'} w_{u,v,k'k} \quad \text{with } \begin{cases} i' = i \times s_h + u \\ j' = j \times s_w + v \end{cases}$$

- $z_{i,j,k}$  is the output of the neuron located in row  $i$ , column  $j$  in feature map  $k$  of the convolutional layer (layer  $l$ ).
- As explained earlier,  $s_h$  and  $s_w$  are the vertical and horizontal strides,  $f_h$  and  $f_w$  are the height and width of the receptive field, and  $f_n$  is the number of feature maps in the previous layer (layer  $l-1$ ).
- $x_{i',j',k'}$  is the output of the neuron located in layer  $l-1$ , row  $i'$ , column  $j'$ , feature map  $k'$  (or channel  $k'$  if the previous layer is the input layer).
- $b_k$  is the bias term for feature map  $k$  (in layer  $l$ ). You can think of it as a knob that tweaks the overall brightness of the feature map  $k$ .
- $w_{u,v,k'k}$  is the connection weight between any neuron in feature map  $k$  of the layer  $l$  and its input located at row  $u$ , column  $v$  (relative to the neuron's receptive field), and feature map  $k'$ .

from the book Hands on Machine Learning with Scikit-Learn and TensorFlow, © Géron, 2017

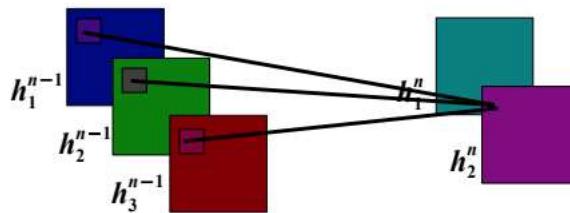
©2018 Peter V. Henstock and David Dowey

## Stacking multiple feature maps

- Repeat about the red boxes once again, it's important...

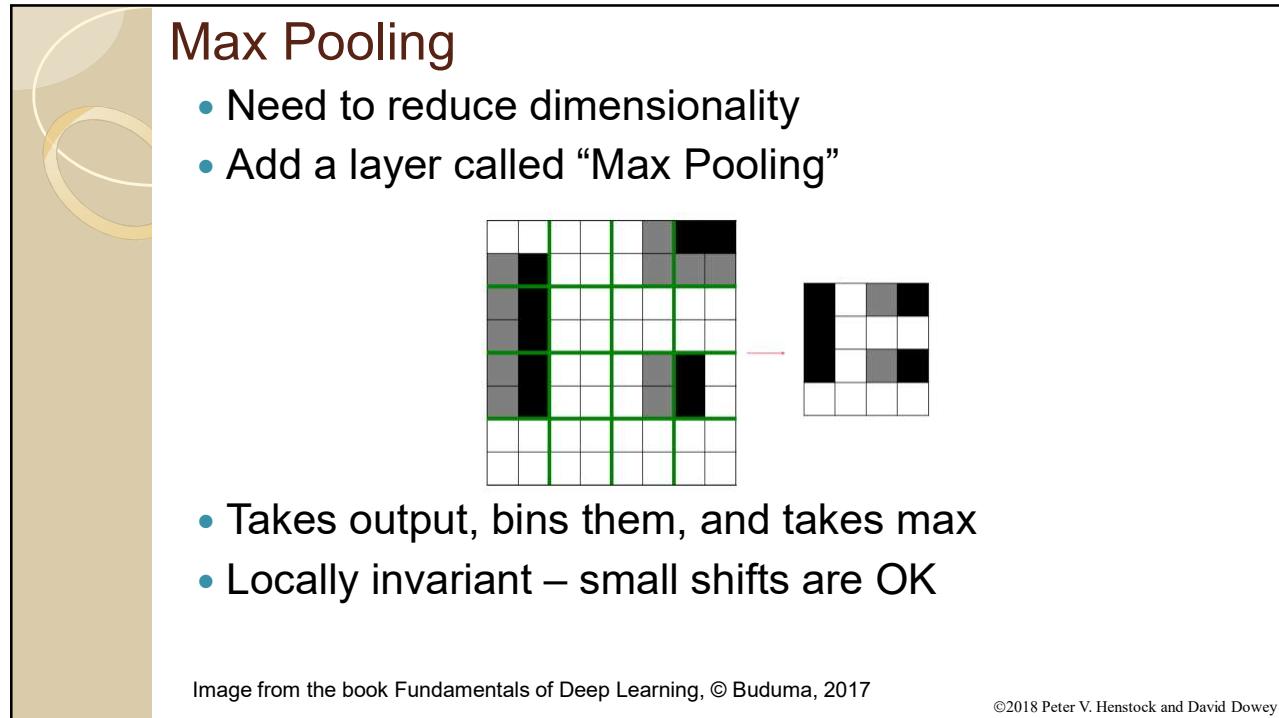
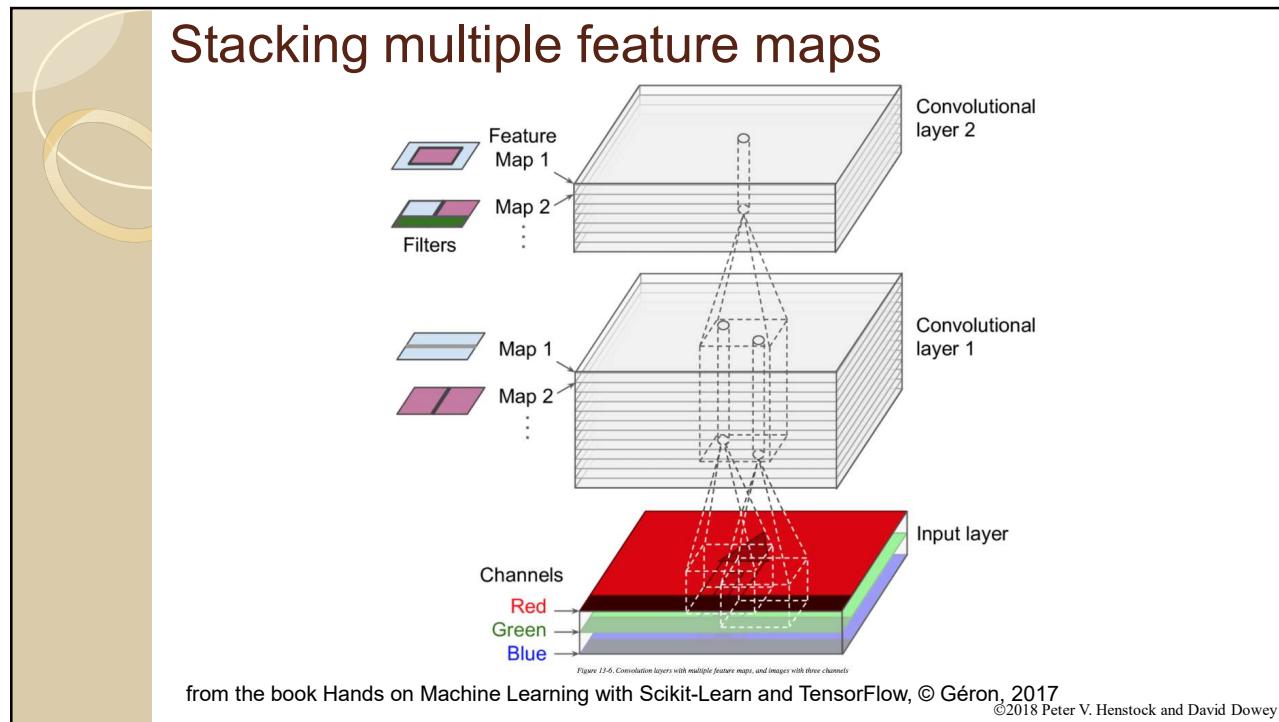
$$h_j^n = \max(0, \sum_{k=1}^K h_k^{n-1} * w_{kj}^n)$$

output  
feature map     
 input feature  
map     
 kernel



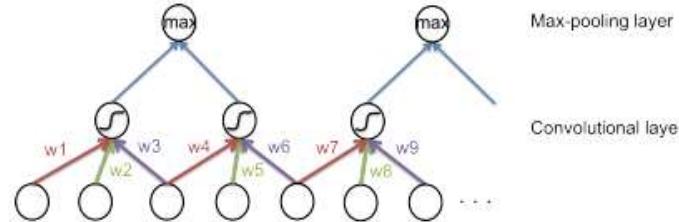
slide 64 of this [CVPR 2014 tutorial](#) by [Marc'Aurelio Ranzato](#)

©2018 Peter V. Henstock and David Dowey



## Max Pooling

Tutorial on Deep Learning Part 2: Autoencoder, Convolutional Neural Networks and Recurrent neural Networks by Quoc V. Le @ Google 2015



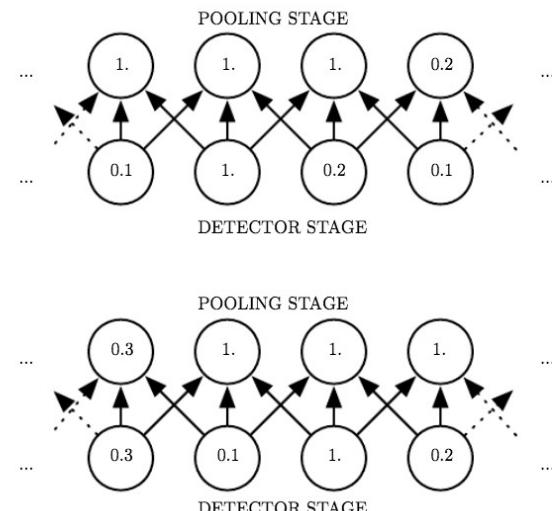
- Interweave max-pooling layer between convolutional layers
- Reduces the scale
- Shift/Translation invariance
  - Events in left or middle → same node
  - Defines local region

©2018 Peter V. Henstock and David Dowey

## Max-Pooling

Introduces  
Shift/Translation  
invariance

- Bottom row is output of activation function
- Unit that fires is most important (1.)
- (right) effect of a 1 position shift is reduced on top row

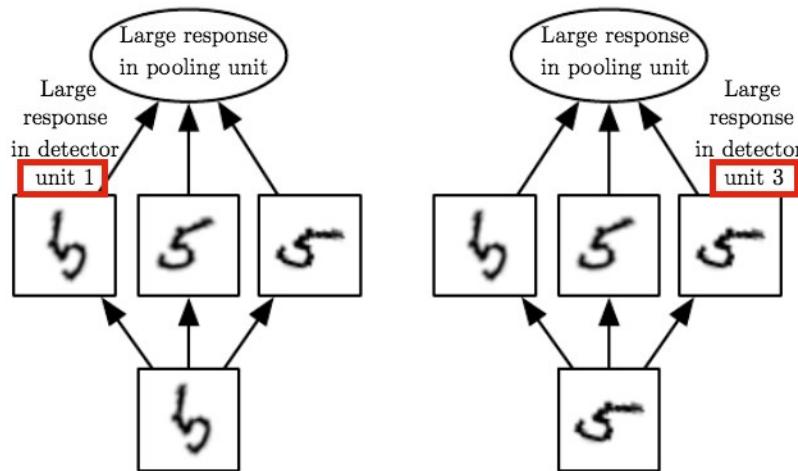


from the book Deep Learning, © Goodfellow et al., 2016

©2018 Peter V. Henstock and David Dowey

## Other Pooling: across multiple channels

Pooling over multiple feature maps with separate learned parameters:  
units learn with transformation invariance: fires on 5 regardless of angle



from the book Deep Learning, © Goodfellow et al., 2016

©2018 Peter V. Henstock and David Dowey

## CNN architectures: Stacked layers of functions

Choices i.e. hyperparameters:

- for layers: width, height, depth
- for filters/kernels: affect learned features
  - spatial extent ( $w, h$ ), depth, stride, zero-padding
- pooling: max-pooling or something else
- activation
- regularization
- initialization
- batch-normalization
- optimizer

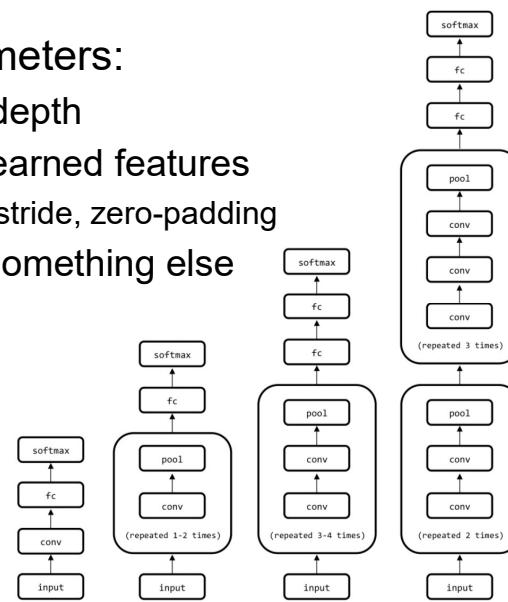
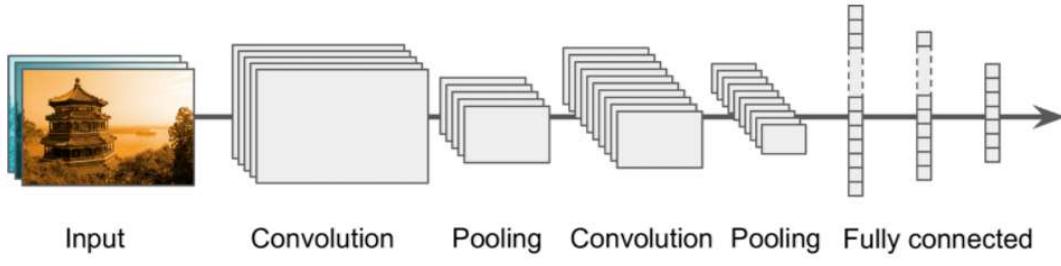


Image from the book Fundamentals of Deep Learning, © Buduma, 2017

©2018 Peter V. Henstock and David Dowey

## CNN architectures: Stacked layers of functions



- Convolutions – filter dimensions, stride
- Multiple shapes, sequences
- Max pooling steps
- Fully-connected layers

from the book Hands on Machine Learning with Scikit-Learn and TensorFlow, © Géron, 2017

©2018 Peter V. Henstock and David Dowey

# Deep Neural Network Details for Implementation

©2018 Peter V. Henstock and David Dowey

## Training CNN deep learning networks

- Initializing the weights – Xavier & He initialization
- Batch normalization

©2018 Peter V. Henstock and David Dowey

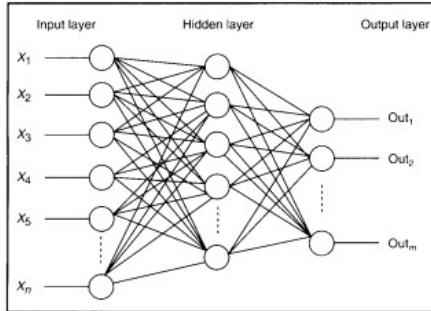
## Backpropagation Algorithm

- Initialize weights random & small—why?

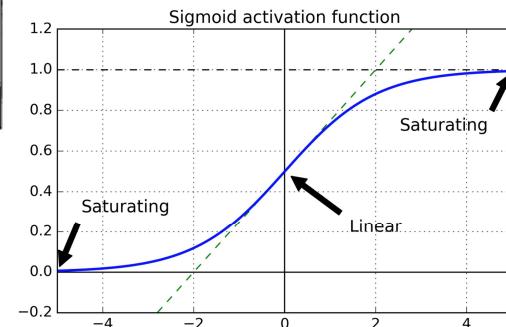
©2018 Peter V. Henstock and David Dowey

## Backpropagation Algorithm

- Initialize weights random & small—why?



- $N(0, 1/\sqrt{\#\text{inputs}})$
- “Symmetry Breaking”



©2018 Peter V. Henstock and David Dowey

## Vanishing and exploding gradients

- Often gradients get smaller and smaller as backprop progresses to the lower layers
- Glorot & Bengio (2010) – major suspects:
  - Sigmoid activation (remember its mean = 0.5)
  - Weight initialisation - How do you think this was being done?
  - Combination of activation function and initialization meant variance outputs each layer > variance inputs each layer

©2018 Peter V. Henstock and David Dowey

## Xavier & He initialization

- Equalize variance of outputs of each layer to its inputs
- Equalize variance of gradients after going through and coming back through layer
- How? --> random initialization of weights using:

Table 11-1. Initialization parameters for each type of activation function

Activation function	Uniform distribution [-r, r]	Normal distribution
Logistic	$r = \sqrt{\frac{6}{n_{\text{inputs}} + n_{\text{outputs}}}}$	$\sigma = \sqrt{\frac{2}{n_{\text{inputs}} + n_{\text{outputs}}}}$
Hyperbolic tangent	$r = 4\sqrt{\frac{6}{n_{\text{inputs}} + n_{\text{outputs}}}}$	$\sigma = 4\sqrt{\frac{2}{n_{\text{inputs}} + n_{\text{outputs}}}}$
ReLU (and its variants)	$r = \sqrt{2}\sqrt{\frac{6}{n_{\text{inputs}} + n_{\text{outputs}}}}$	$\sigma = \sqrt{2}\sqrt{\frac{2}{n_{\text{inputs}} + n_{\text{outputs}}}}$

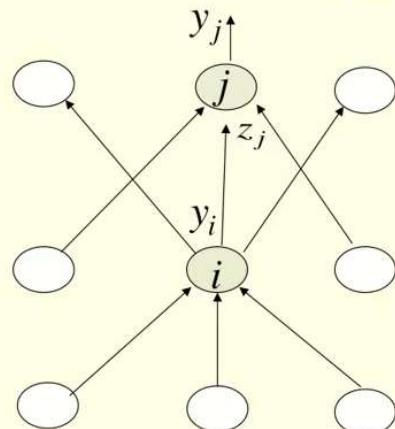
- Speeds up training of deep networks

from the book Hands on Machine Learning with Scikit-Learn and TensorFlow, © Géron, 2017

©2018 Peter V. Henstock and David Dowey

## Backpropagation on one slide

### Backpropagating $dE/dy$



$$\frac{\partial E}{\partial z_j} = \frac{dy_j}{dz_j} \frac{\partial E}{\partial y_j} = y_j (1 - y_j) \frac{\partial E}{\partial y_j}$$

$$\frac{\partial E}{\partial y_i} = \sum_j \frac{dz_j}{dy_i} \frac{\partial E}{\partial z_j} = \sum_j w_{ij} \frac{\partial E}{\partial z_j}$$

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial z_j}{\partial w_{ij}} \frac{\partial E}{\partial z_j} = y_i \frac{\partial E}{\partial z_j}$$

- Just the chain rule
- Hinton, Srivastava and Swersky slides

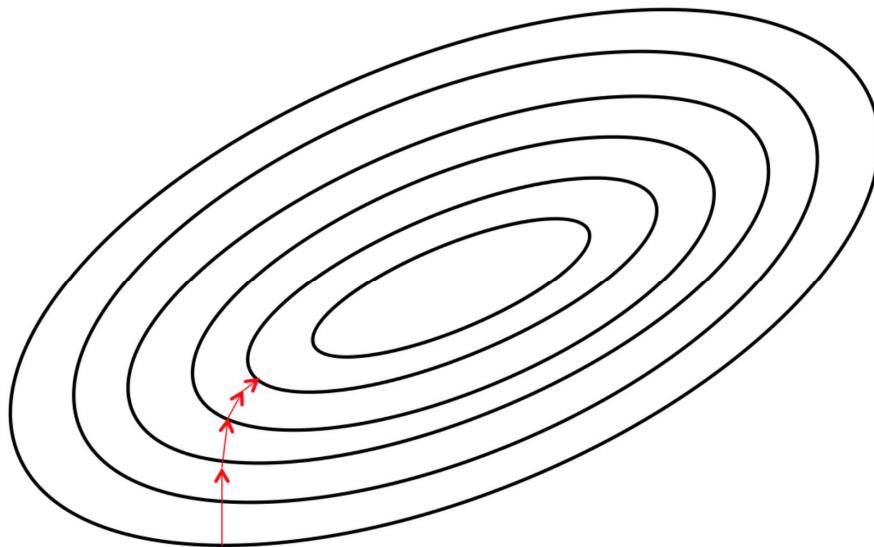
©2018 Peter V. Henstock and David Dowey

## Backpropagation Algorithm

- Initialize weights random & small
- Loop
  - For each training example
    - Put into network
    - Compute network outputs
    - For each output unit  $k$ 
      - $d_k = y_k(1-y_k)(t_k - y_k)$
    - For each hidden unit
      - $d_h = o_h(1-o_h) \sum_{k=\text{out}} d_k w_{hk}$  where  $o_h$  = output of hidden
    - Update weights  $w_{ij} \leftarrow w_{ij} + \Delta\alpha d_j x_{ij}$

©2018 Peter V. Henstock and David Dowey

## Gradient Descent



from the book Fundamentals of Deep Learning, © Buduma, 2017

©2018 Peter V. Henstock and David Dowey

## Gradient Descent – it's slow

Why?

©2018 Peter V. Henstock and David Dowey

## Gradient Descent – it's slow

- 1) Need to perform gradient descent over the entire network of weights
  - Multiple layers x wide inputs in many cases
    - 100x100 images → 10K weights from input
- 2) Not guaranteed convexity (global min)
  - What to do about this?
- 3) Sensitive: use slow learning rate
- 4) Present inputs & outputs sequentially
  - Epoch = one round through data
  - May need 1000s of epochs

©2018 Peter V. Henstock and David Dowey

## Batch Mode

- Initialize all the weights  $w_i$
- Loop
  - Present each training instance to the network
    - Input  $x$  vector
    - Compute the output  $y$
  - Compare  $y$  to the  $t$  from the training data
  - Compute  $\Delta w_i = \sum_{k=inputs} (t_k - y_k) (-x_{ki})$
  - Update the new  $w_i$  weights
- Epoch = presentation of all training data
- Batch mode updates after each epoch

©2018 Peter V. Henstock and David Dowey

## How to Speed Gradient Descent?

- Idea: vary amount of data used to compute the perfect gradients
- 1) Batch gradient descent
    - Use epoch of data to update weights
    - Can't update with new data coming in
  - 2) Stochastic gradient descent
    - Update weights with each training sample
    - Wrong & overshoots, but OK for small steps
  - 3) Mini-batch
    - Update weights every  $n$  samples    $n \sim 256$

©2018 Peter V. Henstock and David Dowey

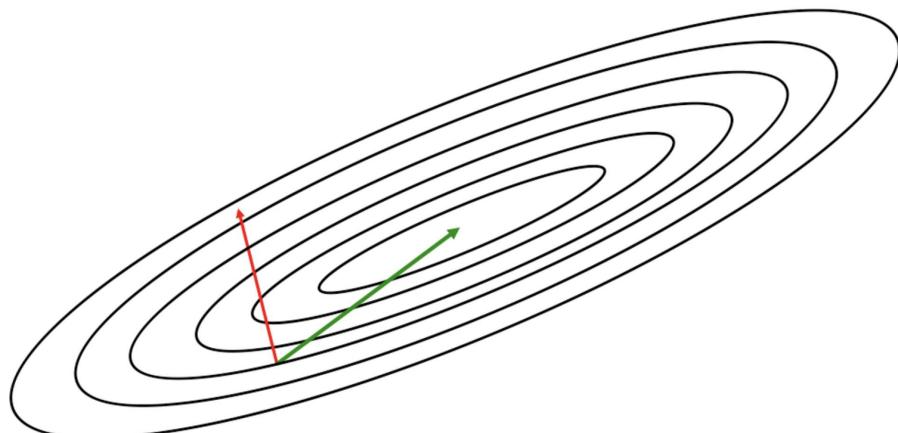
## Batch vs. Incremental Grad. Desc.

- Batch
  - + Good for well behaved error surfaces
  - - Slow to update
  - + Can get stuck in saddle-points (flat parts) and local minima
- Incremental (stochastic, "online", but also mini-batch)
  - + Faster convergence
  - + Can help avoid local minima
  - - Lose the guarantees of convergence

©2018 Peter V. Henstock and David Dowey

## Transforming Inputs

- Scale the inputs so they are fairly equivalent in scale



- <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.159.7105&rep=rep1&type=pdf>
- [https://en.wikipedia.org/wiki/Gradient\\_descent](https://en.wikipedia.org/wiki/Gradient_descent)

from the book Fundamentals of Deep Learning, © Buduma, 2017

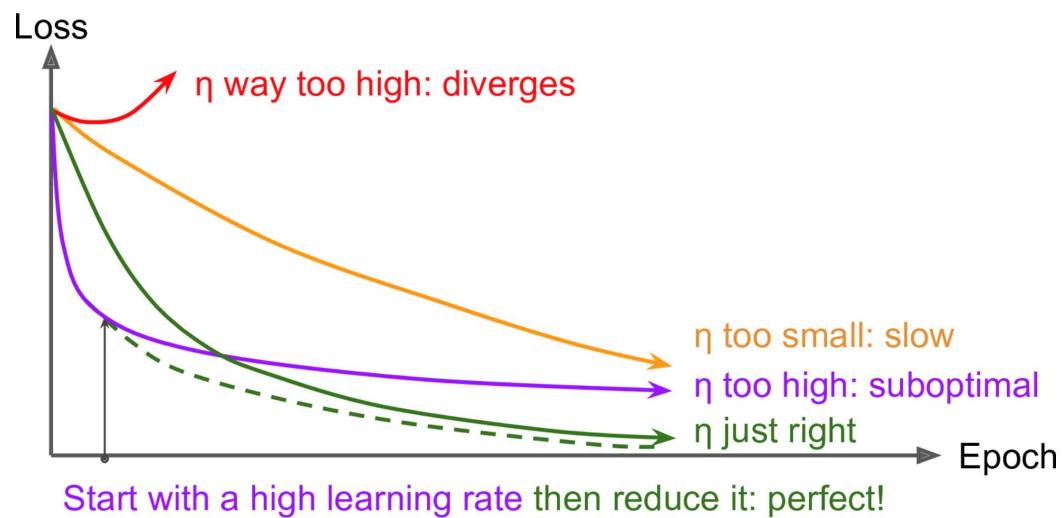
©2018 Peter V. Henstock and David Dowey

## How to Speed Gradient Descent?



©2018 Peter V. Henstock and David Dowey

## How to Speed Gradient Descent?



from the book Hands on Machine Learning with Scikit-Learn and TensorFlow, © Géron, 2017  
©2018 Peter V. Henstock and David Dowey

## Momentum

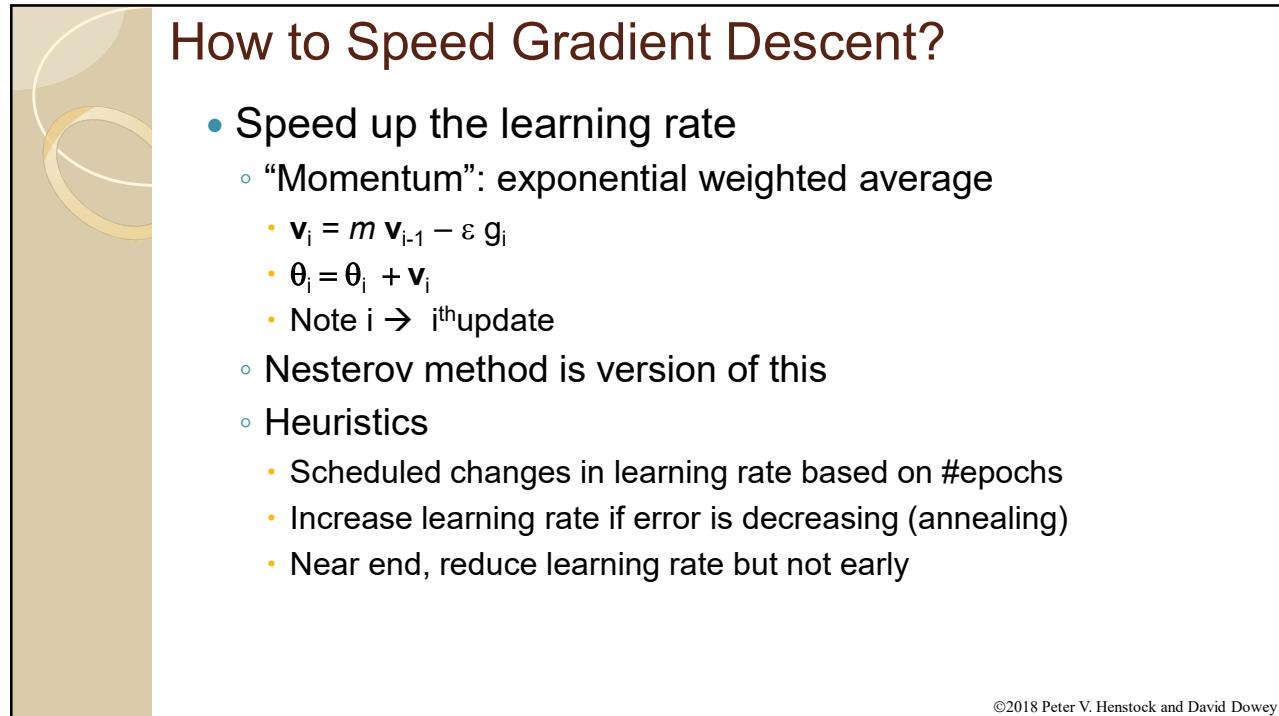
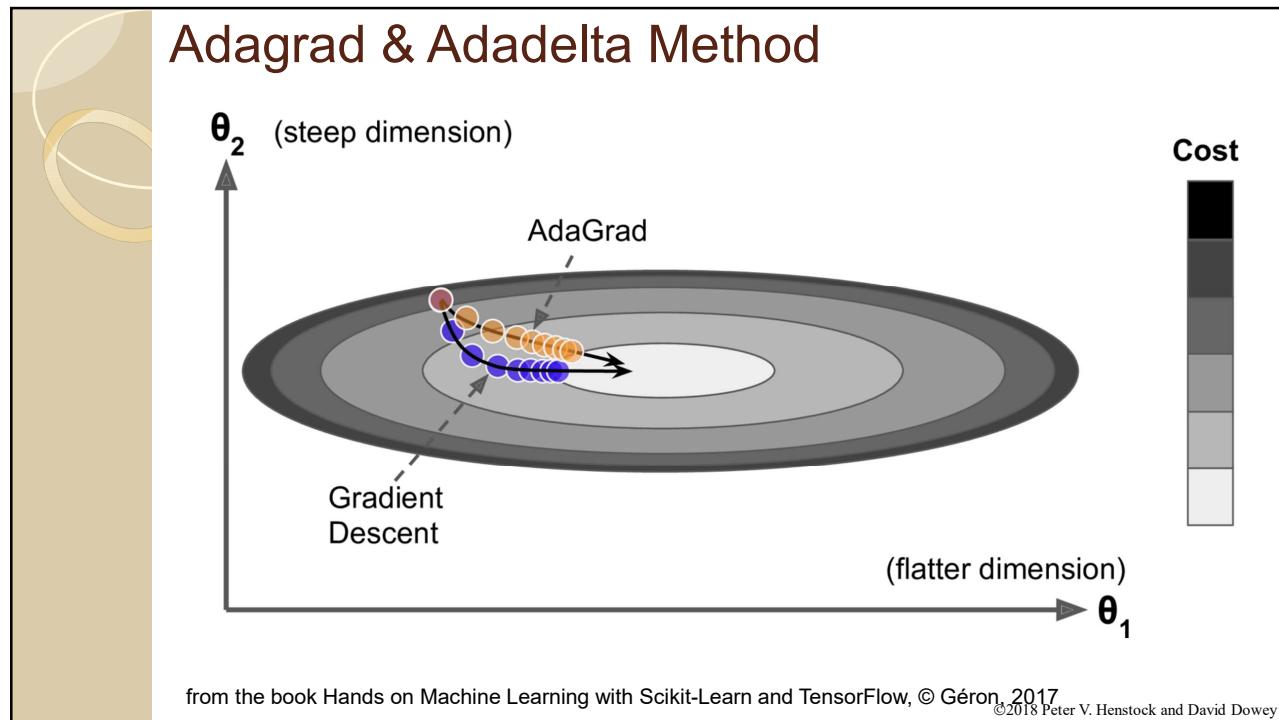
- Simple extension of SGD – used for many years
- Helps in the case of a long, narrow valley (SGD slow)
  - Accelerate along dimensions in which gradient consistently points in the same direction
  - Slow down along dimensions where the sign of the gradient continues to change
- How? Keep track of past parameter updates with exponential decay

©2018 Peter V. Henstock and David Dowey

## Adagrad & Adadelta Method

- Adagrad
  - Provides a gradient descent helper
  - Adjusts the step size based on the slope:
  - Large gradients = small learning rate & vice versa
  - Uses the stochastic updating
  - Applies different learning rate for each parameter separately
  - Incorporates past gradients as well
  - Accumulation with sum of squared past gradients:
    - learning rate continues to shrink - can → 0 change
- Adadelta
  - Fix to Adagrad using window of past gradients that fixes the → 0 issue

©2018 Peter V. Henstock and David Dowey



## Gradient Descent – it's slow

Why?

©2018 Peter V. Henstock and David Dowey

## How to Speed Gradient Descent?

- Momentum optimization,
- Nesterov Accelerated Gradient,
- AdaGrad & Adadelta
- RMSProp
  - Decay parameter, like Adagrad accumulate squared past gradients
- Adam optimization.
  - Momentum and RMSProp combined

©2018 Peter V. Henstock and David Dowey

## How to Speed Gradient Descent?

*Equation 11-8. Adam algorithm*

1.  $\mathbf{m} \leftarrow \beta_1 \mathbf{m} - (1 - \beta_1) \nabla_{\theta} J(\theta)$
2.  $\mathbf{s} \leftarrow \beta_2 \mathbf{s} + (1 - \beta_2) \nabla_{\theta} J(\theta) \otimes \nabla_{\theta} J(\theta)$
3.  $\hat{\mathbf{m}} \leftarrow \frac{\mathbf{m}}{1 - \beta_1^t}$  This only for information purposes
4.  $\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \beta_2^t}$  3 and 4 are correction hyperparameters to account for initialization bias
5.  $\theta \leftarrow \theta + \eta \hat{\mathbf{m}} \oslash \sqrt{\hat{\mathbf{s}}} + \epsilon$

•  $t$  represents the iteration number (starting at 1).

from the book Hands on Machine Learning with Scikit-Learn and TensorFlow, © Géron, 2017  
 ©2018 Peter V. Henstock and David Dowey

## Batch Normalization

- What if the outputs of the neurons in the bottom part of the network shift?

©2018 Peter V. Henstock and David Dowey

## Batch Normalization

- What if the outputs of the neurons in the bottom part of the network shift?
- What would happen to the upper layers?

©2018 Peter V. Henstock and David Dowey

## Batch Normalization

shifts in the distribution of output from the neurons can have impact all the way up from the bottom-most to the top – like stacked boxes

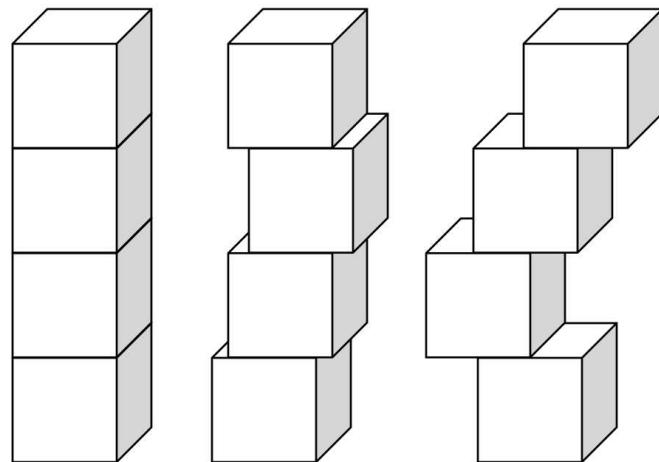
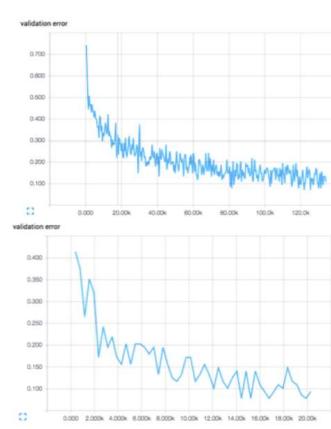


Figure 5-15. When blocks in a layer become shifted too drastically so that they no longer align, the structure can become very unstable

from the book Fundamentals of Deep Learning, © Buduma, 2017

©2018 Peter V. Henstock and David Dowey

## Batch Normalization

- What if the outputs of the neurons in the bottom part of the network shift?
- What would happen to the higher layers?
  - Would have to remember old data
  - Would have to accommodate new data
- How could you solve the problem?

©2018 Peter V. Henstock and David Dowey

## Batch Normalization

- Take logits out of layer before go into the non-linearity
- Z-Normalize each tuple of vector across all examples of mini-batch
- May have to modify the momentum calculations
- Claims:
  - Avoids need for regularization & dropout
  - Several fold faster learning

©2018 Peter V. Henstock and David Dowey

## Overfitting

- Major problem in neural networks
- Lots of parameters to tweak
- How would you avoid it?

©2018 Peter V. Henstock and David Dowey

## Avoid Overfitting

### Regularization: penalize large weights

- Add an extra loss term for the overall model loss function to optimize
  - $L_1$  or  $L_2$
  - "Weight-decay"
  - $Error + \lambda * f(Weights) \rightarrow$  each weight  $\frac{1}{2}\lambda w^2$
  - Need to modify the derived weighting
  - Widely used
- "Max-norm": constrain incoming weights at each neuron by "clipping" them such that the  $L_2$  norm is less than  $r$  - in gradient descent, project vector back

©2018 Peter V. Henstock and David Dowey

## Avoid Overfitting

### Share weights

- We are already doing that as part of CNN

©2018 Peter V. Henstock and David Dowey

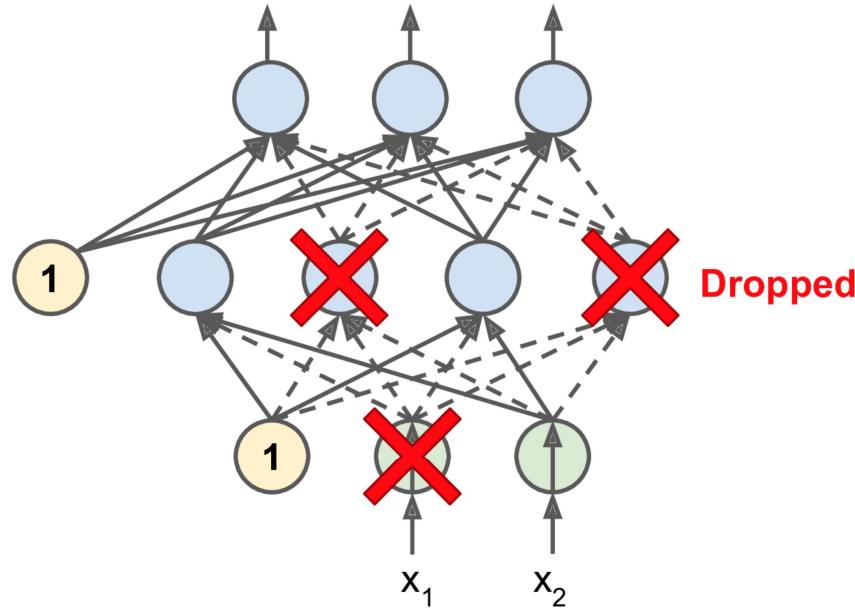
## Avoid Overfitting

### Dropout: randomly nodes & connections during training

- Randomly remove half of hidden nodes
  - Actually set the weights to 0
- Train network normally with updates
- Restore the nodes for next mini-batch

©2018 Peter V. Henstock and David Dowey

## Avoid Overfitting - Dropout



from the book Hands on Machine Learning with Scikit-Learn and TensorFlow, © Géron, 2017  
©2018 Peter V. Henstock and David Dowey

## Avoid Overfitting - Dropout

### Dropout (cont.):

- Randomly remove half of hidden nodes
  - Actually set the weights to 0
  - But if you set half the weights to 0, what will happen to the subsequent activation?
- How to fix it?

©2018 Peter V. Henstock and David Dowey

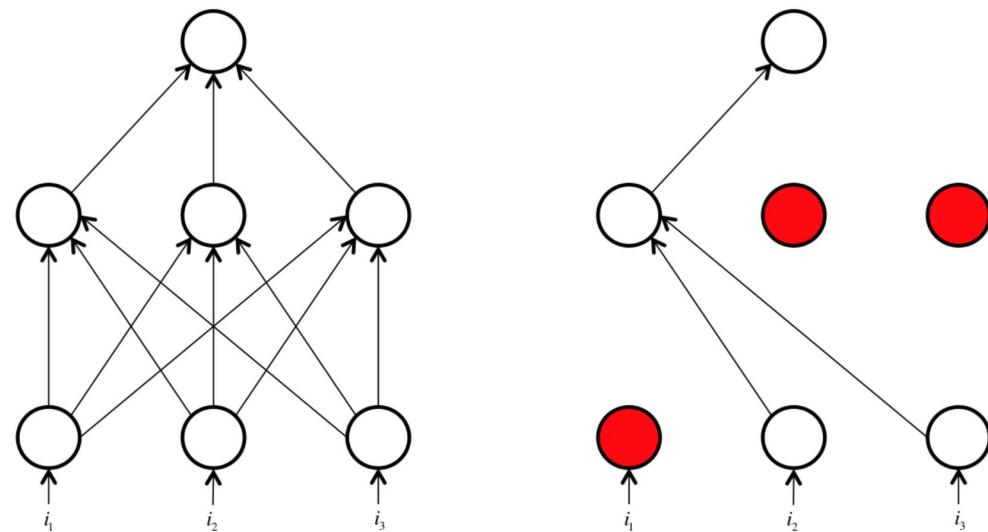
## Avoid Overfitting - Dropout

Dropout (cont.):

- Randomly remove half of hidden nodes
  - Actually set the weights to 0
  - But if you set half the weights to 0, what will happen to the subsequent activation?
- How to fix it?
  - When you run the trained model on the test data, multiply each neuron's input weight by the "keep probability" ( $1-p$ )
  - Or, during training, divide the output weights by the "keep probability" [preferred]

©2018 Peter V. Henstock and David Dowey

## Avoid Overfitting - Dropout



from the book Fundamentals of Deep Learning, © Buduma, 2017

©2018 Peter V. Henstock and David Dowey

## Avoid Overfitting

Stop early by tracking error as  $f(\text{epoch})$

- ? Use model averaging (like bagging)
- ? Use cross-validation to figure out when training error increases

©2018 Peter V. Henstock and David Dowey

## Hyperparameter tuning:

- How do you know what is a good hyperparameter choice?

©2018 Peter V. Henstock and David Dowey

## Hyperparameter tuning:

- How do you know what is a good hyperparameter choice?
- Capacity: ability to fit a wide variety of functions

©2018 Peter V. Henstock and David Dowey

## Hyperparameter tuning:

- Training error – lower is better
- Generalization error – even lower is better, and you don't want this to start increasing
- Resource constraints – computing power, runtime, memory

©2018 Peter V. Henstock and David Dowey

## Hyperparameter tuning:

- training error – lower is better
- generalization error – even lower is better
- resource constraints – computing power, runtime, memory

The effective “capacity” of a model depends on:

- representational capacity of model (e.g. # layers, depth)
- success of algorithm to minimize cost function in training
  - the ability to discover the functions that best give representation
- regularization of the cost function and training procedure

©2018 Peter V. Henstock and David Dowey

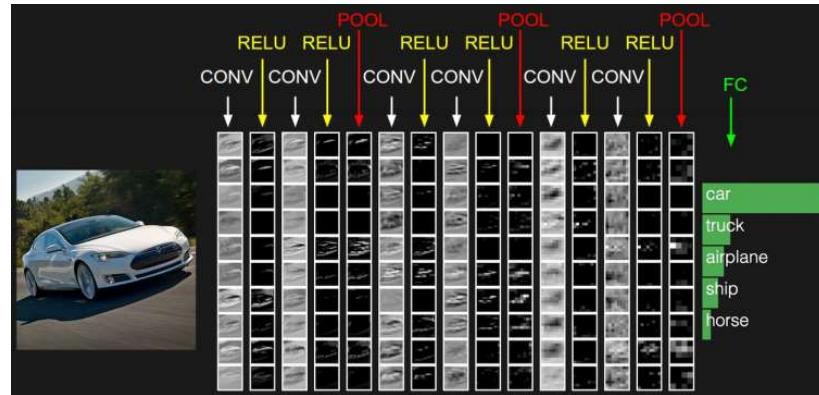
Hyper-parameter	Increases capacity when	Reason	Caveats
Number of hidden units	increased	Increasing the number of hidden units increases the representational capacity of the model.	Increasing the number of hidden units increases both the time and memory cost of essentially every operation on the model.
Learning rate	tuned optimally	An improper learning rate, whether too high or too low, results in a model with low effective capacity due to optimization failure.	
Convolution kernel width (i.e. spatial extent, $w \times h$ )	increased	Increasing the kernel width increases the receptive field and therefore the number of parameters in the model = more distinctive features	A wider kernel results in a narrower output dimension, reducing model capacity unless you use implicit zero padding to reduce this effect. Wider kernels require more memory for parameter storage and increase runtime, but a narrower output reduces memory cost
Implicit zero padding	increased	Adding implicit zeros before convolution keeps the representation size large.	Increases time and memory cost of most operations.
Weight decay coefficient	decreased	Decreasing the weight decay coefficient frees the model parameters to become larger	
Dropout rate	decreased	Dropping units less often gives the units more opportunities to “conspire” with each other to fit the training set.	

from the book Deep Learning p.426, © Goodfellow et al., 2016

©2018 Peter V. Henstock and David Dowey

## Putting it all together

- <http://cs231n.github.io/convolutional-networks/>



©2018 Peter V. Henstock and David Dowey

# Deep Learning Toolsets & Architectures

©2018 Peter V. Henstock and David Dowey

## Existing DeepLearning Tools

- Theano (U Montreal) more for research
- Torch (Facebook) written in Lua
- Caffe (Berkeley)
- Neon
- TensorFlow (Google)
- CNTK (Microsoft)
- MXNet (Apache)
- Front-ends
  - Keras (Tensorflow , CNT, Theano)
  - PyTorch (Torch)

©2018 Peter V. Henstock and David Dowey

## TensorFlow

- Released in 2015 by Google – used by Google Brain
- APIs for Python C/C++ Java & Go
- Tensor = multidimensional array
  - How data are stored and enter into computations
- Includes TensorBoard for visualizing learning through the network graphs, scalars and histograms
- Works out-of-the-box on multiple CPUs and GPUs with acceleration (Nvidia CUDA GPUs)
- Keras front-end is now built in: tf.keras
- current release 1.11

©2018 Peter V. Henstock and David Dowey

## Resources

- TensorFlow:
  - <https://www.tensorflow.org/tutorials/>
  - <https://developers.google.com/machine-learning/crash-course/>
  - <https://github.com/tensorflow/models>  
<https://www.tensorflow.org/hub/>
  - <https://keras.io/applications/>
- Kaggle kernels
  - <https://www.kaggle.com/kernels>
- Model repositories, etc.
  - [https://mxnet.apache.org/model\\_zoo/index.html](https://mxnet.apache.org/model_zoo/index.html)
  - <https://modelzoo.co/>
  - <https://modeldepot.io/>

©2018 Peter V. Henstock and David Dowey

## Types of layers in a CNN

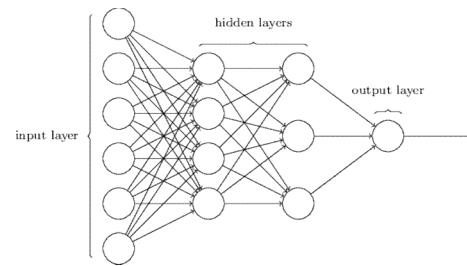
- Convolution
  - Activation, Batch Normalization
- Pooling
- Dropout
- Fully Connected

Remember – what constitutes a "layer"?

©2018 Peter V. Henstock and David Dowey

## Layers

We have gone from this:



to this:

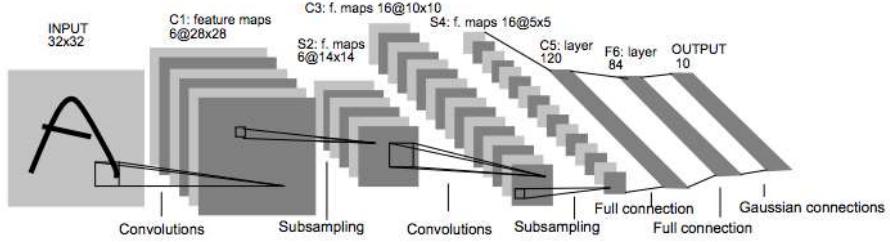


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

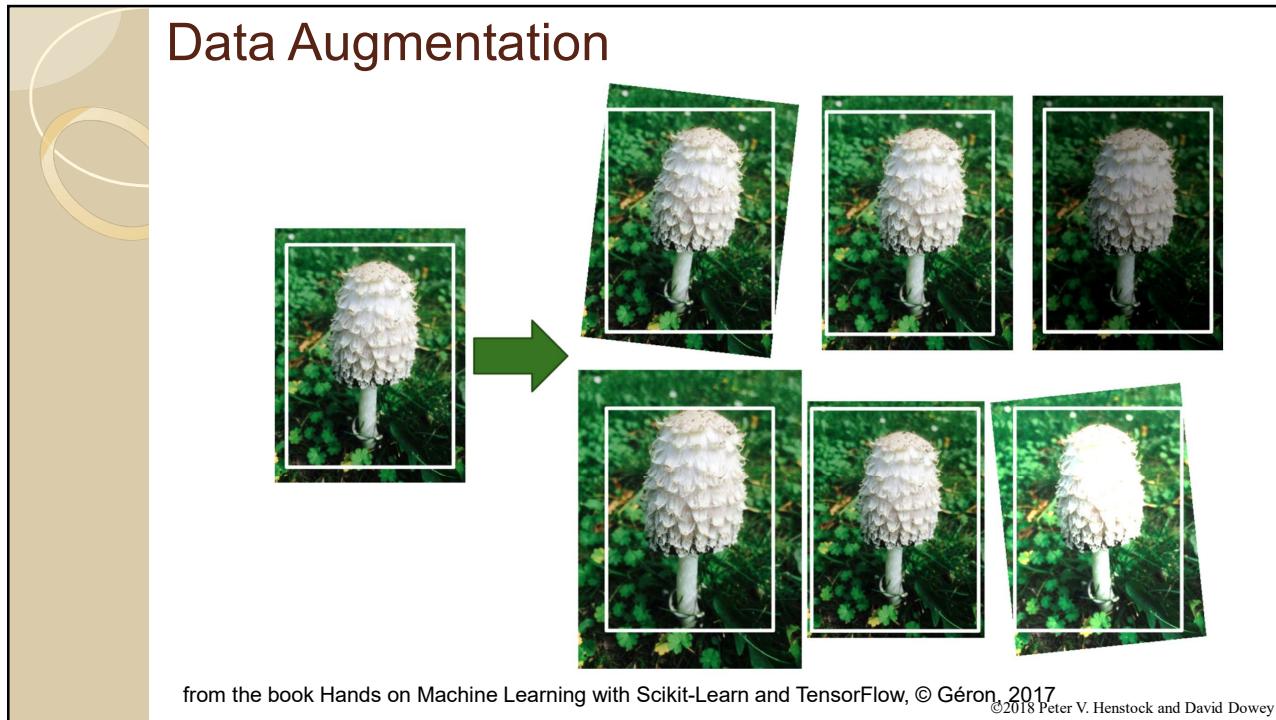
Convention: layers are the ones with weights

©2018 Peter V. Henstock and David Dowey

## CNN preprocessing pipelines

- normalize the images – all images to same scale
  - all to [0,1] or [-1,1] but not mixing [0,1] **and** [0,255]
- most architectures require to have standard size
  - cropping, scaling to fit – dynamic pooling can adjust
- dataset augmentation
  - reduce generalization error by finding more training samples
- bounding box image segmentation into locations

©2018 Peter V. Henstock and David Dowey



## Data Augmentation

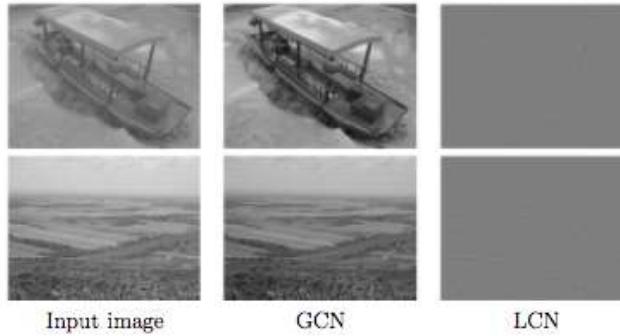
Essential to prevent overfitting in small datasets

- Mirroring – horizontal, vertical
- Rotations
- Random cropping
- Shearing
- Local warping
- Width/height shift
- Color shifting – add distortion to R G B channels
- PCA color augmentation – keep tint the same (AlexNet)

©2018 Peter V. Henstock and David Dowey

## Contrast normalization (bright – dark pixels)

- Global Contrast Normalization 
$$X'_{i,j,k} = s \frac{X_{i,j,k} - \bar{X}}{\max \left\{ \epsilon, \sqrt{\lambda + \frac{1}{3rc} \sum_{i=1}^r \sum_{j=1}^c \sum_{k=1}^3 (X_{i,j,k} - \bar{X})^2} \right\}}$$
- Whitening: all Prin. Components of PCA have equal variance
- Local Contrast Normalization



from the book Deep Learning, © Goodfellow et al., 2016

©2018 Peter V. Henstock and David Dowey

## Some useful CNN architectures

- LeNet-5 – (arguably) the first CNN architecture - digits
- AlexNet – unparalleled performance on ImageNet
- GoogLeNet – Inception module:optimal local ConvNets
- ResNet – add skip connections to learn on residuals

©2018 Peter V. Henstock and David Dowey

## CNN architectures: LeNet-5 1998

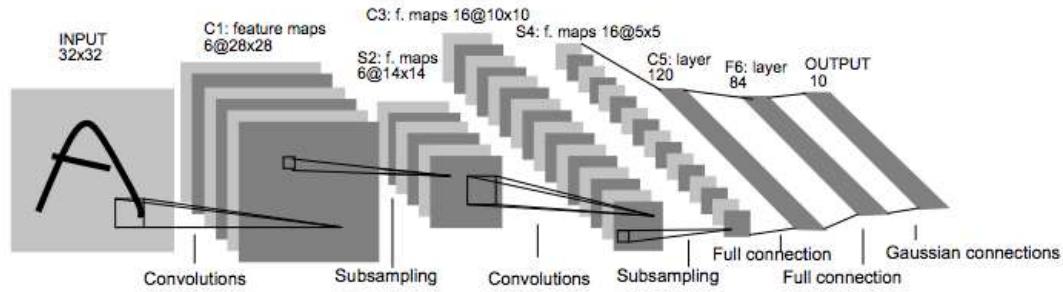


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

LeNet-5 schema from "Gradient-based learning applied to document recognition", LeCun et al, 1998

©2018 Peter V. Henstock and David Dowey

## CNN architectures: LeNet-5 1998

Layer	Activation output		Number of parameters
	Shape	Size	
Input	(32,32,3)	3,072	0
Convolution Layer 1 (f=5, s=1)	(28,28,8)	6,272	208
Pooling Layer 1	(14,14,8)	1,568	0
Convolution Layer 5 (f=5, s=1)	(10,10,16)	1,600	416
Pooling Layer 2	(5,5,16)	40	0
Fully Connected Layer 3	(120,1)	120	48,001
Fully Connected Layer 4	(84,1)	84	10,081
Softmax Layer (output)	(10,1)	10	841

©2018 Peter V. Henstock and David Dowey

## CNN architectures: AlexNet 2012: 17% top-5

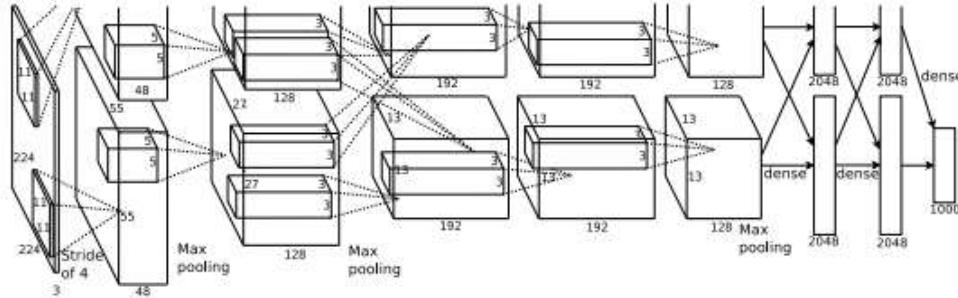


Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input is 150,528-dimensional, and the number of neurons in the network's remaining layers is given by 253,440–186,624–64,896–64,896–43,264–4096–4096–1000.

ImageNet schema from "ImageNet Classification with Deep Convolutional Neural Networks", Krizhevsky et al, 2012

©2018 Peter V. Henstock and David Dowey

## CNN architectures: LeNet-5 and AlexNet

Table 13-1. LeNet-5 architecture

Layer Type	Maps	Size	Kernel size	Stride	Activation
Out	Fully Connected	–	10	–	–
F6	Fully Connected	–	84	–	–
C5	Convolution	120	1 × 1	5 × 5	1
S4	Avg Pooling	16	5 × 5	2 × 2	2
C3	Convolution	16	10 × 10	5 × 5	1
S2	Avg Pooling	6	14 × 14	2 × 2	2
C1	Convolution	6	28 × 28	5 × 5	1
In	Input	1	32 × 32	–	–

Table 13-2. AlexNet architecture

Layer Type	Maps	Size	Kernel size	Stride	Padding	Activation
Out	Fully Connected	–	–	–	–	Softmax
F9	Fully Connected	–	–	–	–	ReLU
F8	Fully Connected	–	–	–	–	ReLU
C7	Convolution	256	13 × 13	3 × 3	1	SAME
C6	Convolution	384	13 × 13	3 × 3	1	SAME
C5	Convolution	384	13 × 13	3 × 3	1	SAME
S4	Max Pooling	256	13 × 13	3 × 3	2	VALID
C3	Convolution	256	27 × 27	5 × 5	1	SAME
S2	Max Pooling	96	27 × 27	3 × 3	2	VALID
C1	Convolution	96	55 × 55	11 × 11	4	SAME
In	Input	3 (RGB)	224 × 224	–	–	–

from the book Hands on Machine Learning with Scikit-Learn and TensorFlow, © Géron, 2017

©2018 Peter V. Henstock and David Dowey

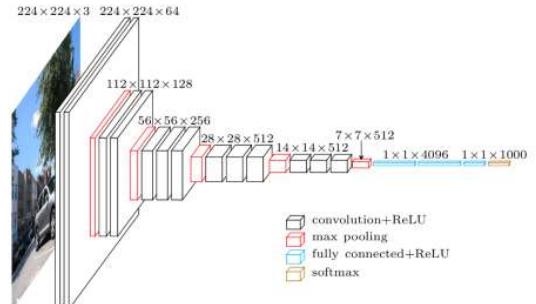
## CNN architectures: VGG-16 VGG-19

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
		maxpool			
conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
		maxpool			
conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256
		maxpool			
conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
		maxpool			
conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
		maxpool			
FC-4096					
FC-4096					
FC-1000					
soft-max					

Table 2: Number of parameters (in millions).

Network	A-A-LRN	B	C	D	E
Number of parameters	133	133	134	138	144

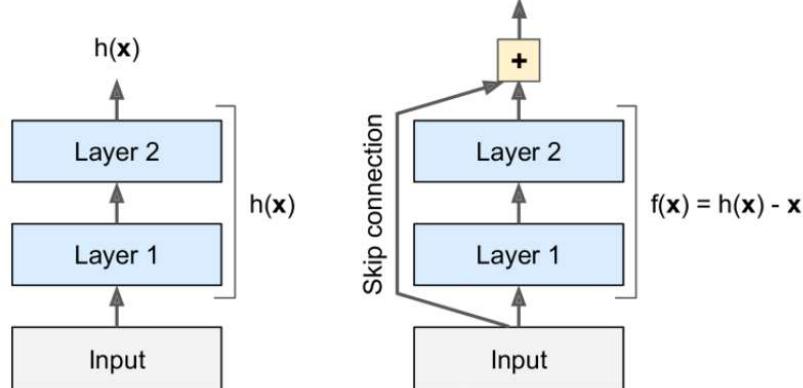
from the paper Simonyan &amp; Zisserman, 2015 VGG - Visual Geometry Group, Oxford



©2018 Peter V. Henstock and David Dowey

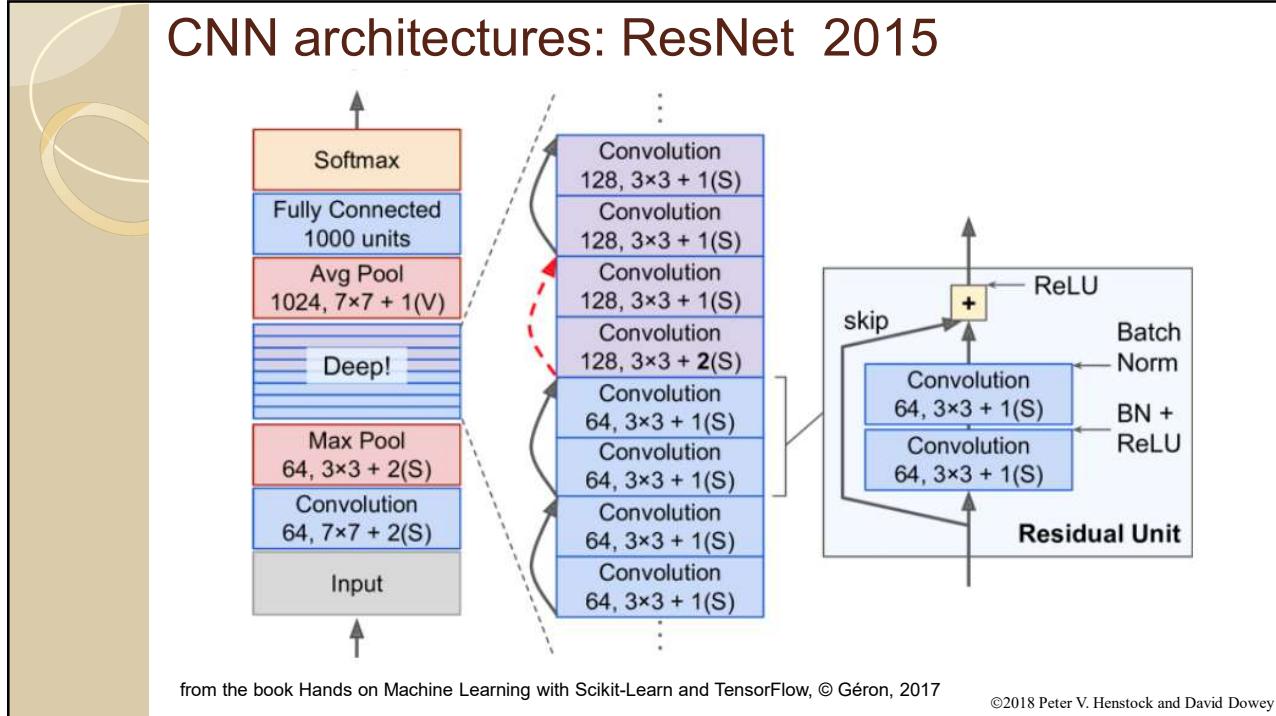
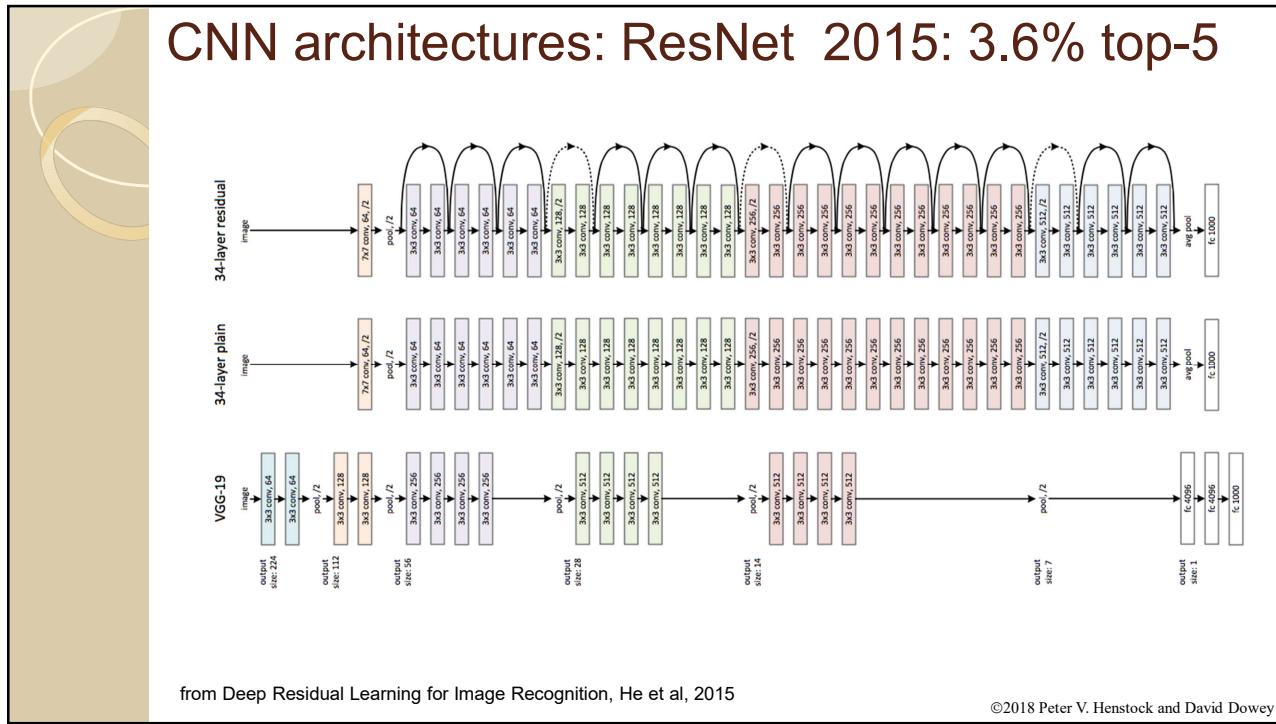
## CNN architectures: ResNet 2015: 3.6% top-5

- Model residual learning through skip connections: add  $x$  to output, network is forced to model  $f(x) = h(x) - x$  instead of  $h(x)$



from the book Hands on Machine Learning with Scikit-Learn and TensorFlow, © Géron, 2017

©2018 Peter V. Henstock and David Dowey



## CNN architectures: ResNet 2015

### Questions

- What happens if the current layer weights are zero (or at least very small)?
- What constraint on the architecture are there on order to make activations "skip" layers?

©2018 Peter V. Henstock and David Dowey

## CNN architectures:

- How do we know how many convolutions, pooling layers, fully-connected layers?
- What size filters works best 5x5, 3x3, 1x1?
- With GPUs and distributed/parallel computing why bother – why not do them all?

©2018 Peter V. Henstock and David Dowey

## CNN architectures: GoogLeNet / Inception: 7%

The diagram illustrates the Inception Module's internal structure. It shows a 'Previous layer' at the bottom, which is processed by three parallel paths. The first path consists of two 1x1 convolutions. The second path has a 1x1 convolution followed by a 3x3 convolution. The third path has a 1x1 convolution followed by a 5x5 convolution. All three paths lead to a '3x3 max pooling' layer. The outputs of all four layers (the three convolutional paths and the max pooling) are concatenated into a 'Filter concatenation' layer at the top.

- "Inception Modules" composed of various sized convolutions
- 1x1 convolutions
- configuration requires patches to align 1x1 3x3 5x5
- Also includes pooling paths
- 22 layers, but
- 12x fewer parameters than AlexNet

Inception Module composition from Going deeper with convolutions, Szegedy et al. 2014  
©2018 Peter V. Henstock and David Dowey

## CNN architectures: GoogLeNet / Inception

### 1 x 1 convolution

- Network-in-network (Lin et al 2012)
- What's the point, anyway?
- How is it different from a fully-connected layer?

The diagram shows the flow of data through a 'Convolutional layer' and a 'CCCP layer'. An 'Input patch (c1 x h x w)' is processed by a 'Convolutional Filter (c2 x c1 x h x w)' to produce an 'Output feature vector (c2 x 1 x 1)'. This is then processed by a 'CCCP layer' using a 'Convolutional Filter (c3 x c2 x 1 x 1)' to produce a final 'Output feature vector (c3 x 1 x 1)'. A callout indicates that the CCCP layer is a 'Representation (Universal Approximator) of the input patch'.

Efficient implementation of CCCP

[http://image-net.org/challenges/LSVRC/2014/slides/ILSVRC2014\\_NUS\\_release.pdf](http://image-net.org/challenges/LSVRC/2014/slides/ILSVRC2014_NUS_release.pdf)  
©2018 Peter V. Henstock and David Dowey

## CNN architectures: GoogLeNet / Inception

### 1 x 1 convolution

- What's the point, anyway?
  - increased network size means better performance
  - there is a limit to how deep you can go → "go deeper"
  - dimension reduction
  - adds non-linearity
- How is it different from a fully-connected layer?
  - at the last layer of a network, it isn't
  - any size input – not the full feature map
    - e.g. LeCun says it gets around fixed input size constraint

©2018 Peter V. Henstock and David Dowey

### 1 x 1 convolutions

E.g. 6 x 6 feature map and single weight, 3, in a 1 x 1 filter with stride 1 : just multiplying every element in the 6 x 6 map by 3. Why bother?

The point: never are doing it with a flat feature map, but a set of feature maps with depth, say, 20 feature maps all stacked in a 6 x 6 x 20 volume.

1 x 1 convolution filter is then 1 x 1 x 20 - with 20 different weights (i.e. numbers) in the filter to be trained, and each one multiplies its own distinct input feature map.

Go from 20 feature maps of size 6 x 6 down to 10 of them.

Each 1 x 1 convolution is like a mini network, since it has weights to learn and an activation function (no pooling or dropout though).

This is dimension reduction - doesn't it make you think of PCA?

©2018 Peter V. Henstock and David Dowey

## 1 x 1 convolutions

Limit to depth of whole network (number of layers)

if you just keep stacking more and more feature maps for each new layer complexity would increase, having more and more weights to learn.

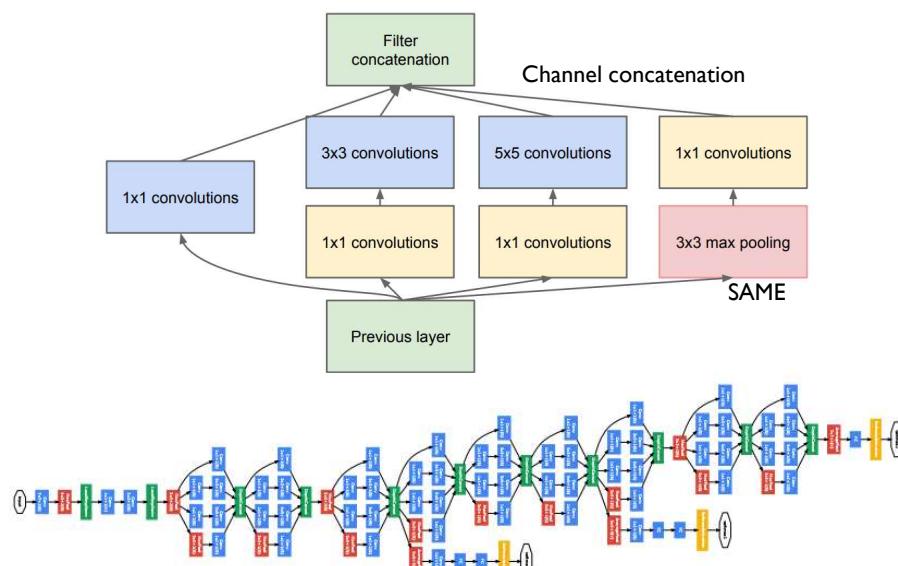
- 1 x 1 convolutions allow you to distill the learning and reduce the depth/complexity,
- they add an additional non-linearity, since they are a mini-network
- not many hyper-parameter choices - you just choose where/when to stick them in, and how many filters to use.
- 1 x 1 convolutions are a lot like flattening the feature maps and doing a Fully Connected layer.

If you used them at the end of a CNN, this is true, they would be essentially the same.

- 1 x 1 convolutions can be used with any sized input, since you are just passing them across all the pixels in the input maps/channels.

©2018 Peter V. Henstock and David Dowey

## CNN architectures: GoogLeNet / Inception



Inception Module composition from Going deeper with convolutions, Szegedy et al. 2014

©2018 Peter V. Henstock and David Dowey

## CNN architectures: GoogLeNet / Inception

type	patch size/ stride	output size	depth	#1×1	#3×3 reduce	#3×3	#5×5 reduce	#5×5	pool proj	params	ops
convolution	7×7/2	112×112×64	1							2.7K	34M
max pool	3×3/2	56×56×64	0								
convolution	3×3/1	56×56×192	2		64	192				112K	360M
max pool	3×3/2	28×28×192	0								
inception (3a)		28×28×256	2	64	96	128	16	32	32	159K	128M
inception (3b)		28×28×480	2	128	128	192	32	96	64	380K	304M
max pool	3×3/2	14×14×480	0								
inception (4a)		14×14×512	2	192	96	208	16	48	64	364K	73M
inception (4b)		14×14×512	2	160	112	224	24	64	64	437K	88M
inception (4c)		14×14×512	2	128	128	256	24	64	64	463K	100M
inception (4d)		14×14×528	2	112	144	288	32	64	64	580K	119M
inception (4e)		14×14×832	2	256	160	320	32	128	128	840K	170M
max pool	3×3/2	7×7×832	0								
inception (5a)		7×7×832	2	256	160	320	32	128	128	1072K	54M
inception (5b)		7×7×1024	2	384	192	384	48	128	128	1388K	71M
avg pool	7×7/1	1×1×1024	0								
dropout (40%)		1×1×1024	0								
linear		1×1×1000	1							1000K	1M
softmax		1×1×1000	0								

Table 1: GoogLeNet incarnation of the Inception architecture

Going deeper with convolutions, Szegedy et al. 2014

©2018 Peter V. Henstock and David Dowey

## Convolutional Neural Networks

- Good

- Top performance for their domain – image recognition
- Representational capacity – highly flexible
- Speed

- Bad

- No memory → LSTM, memory networks, etc.
- Fixed input size: good for images and other fixed-length vectors
- Fixed output size – i.e. categories → extend pre-trained networks, outputs as features
- Every new technique seems to add more hyper-parameters

©2018 Peter V. Henstock and David Dowey

## Convolutional Neural Networks

- **Ugly**

- "Just made it up" – how to choose architecture?
  - Activation functions are at the heart but we just use what works
  - *"One must be cautious though: although the proposed architecture has become a success for computer vision, it is still questionable whether its quality can be attributed to the guiding principles that have lead to its construction."* (Going deeper with convolutions)
- Throwing away information: Is pooling needed?
  - equivariance not invariance
  - reflect the transformation as information
  - disentangle not discard (Hinton)

©2018 Peter V. Henstock and David Dowey

## Transfer Learning

- Not everyone can train an Inception-v3 CNN with ImageNet images for 3 weeks with multiple GPUs
- Three scenarios:
  1. Use a pre-trained CNN → move on from checkpoint
  2. Fine tune the CNN: replace the classifier at the top of the CNN and continue to train – backpropagation will fine tune the weights; keep some layers frozen
  3. Use the CNN as a feature generator
    - Take a trained CNN, remove the last FC layer and use the output layer as features e.g. AlexNet = 4096 features
    - Train a classical linear classifier for the new images

©2018 Peter V. Henstock and David Dowey

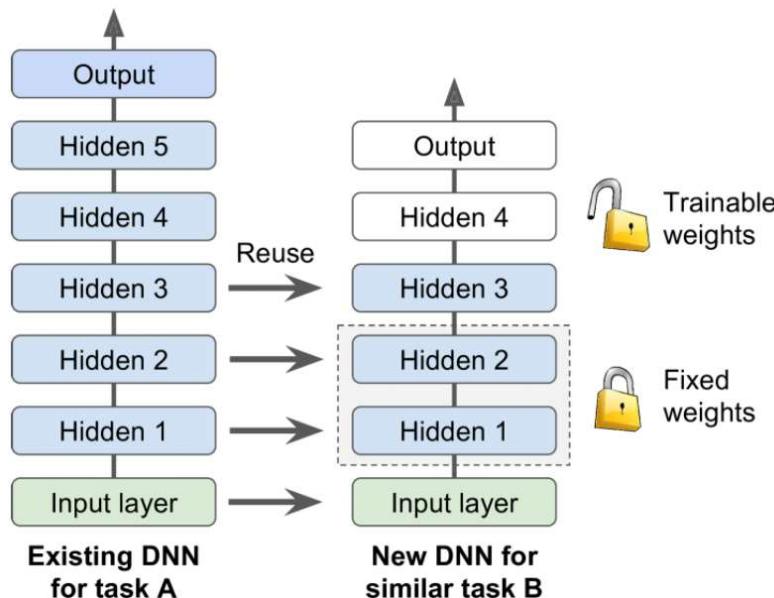
## Transfer Learning

Important question: how much data do you have?

- Lots: Options 1 and 2 – pre-train & move on or fine tune the whole, or parts of, the base
- Not much: Option 3 – use a pre-trained CNN as a feature generator
  - 1. Run base once - save output to disk – use as input layer
  - 2. Re-run the (frozen) base every time – more computations, but allows you to do data augmentation
- Another important question: do you have a GPU?

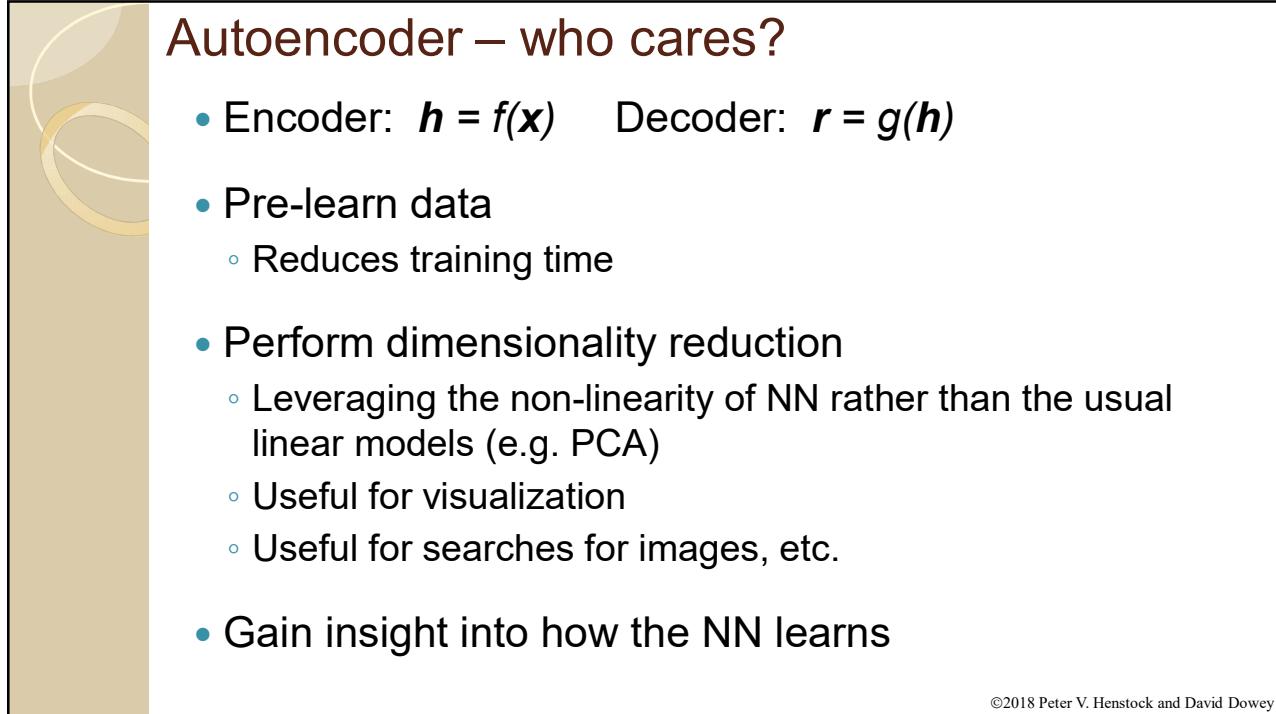
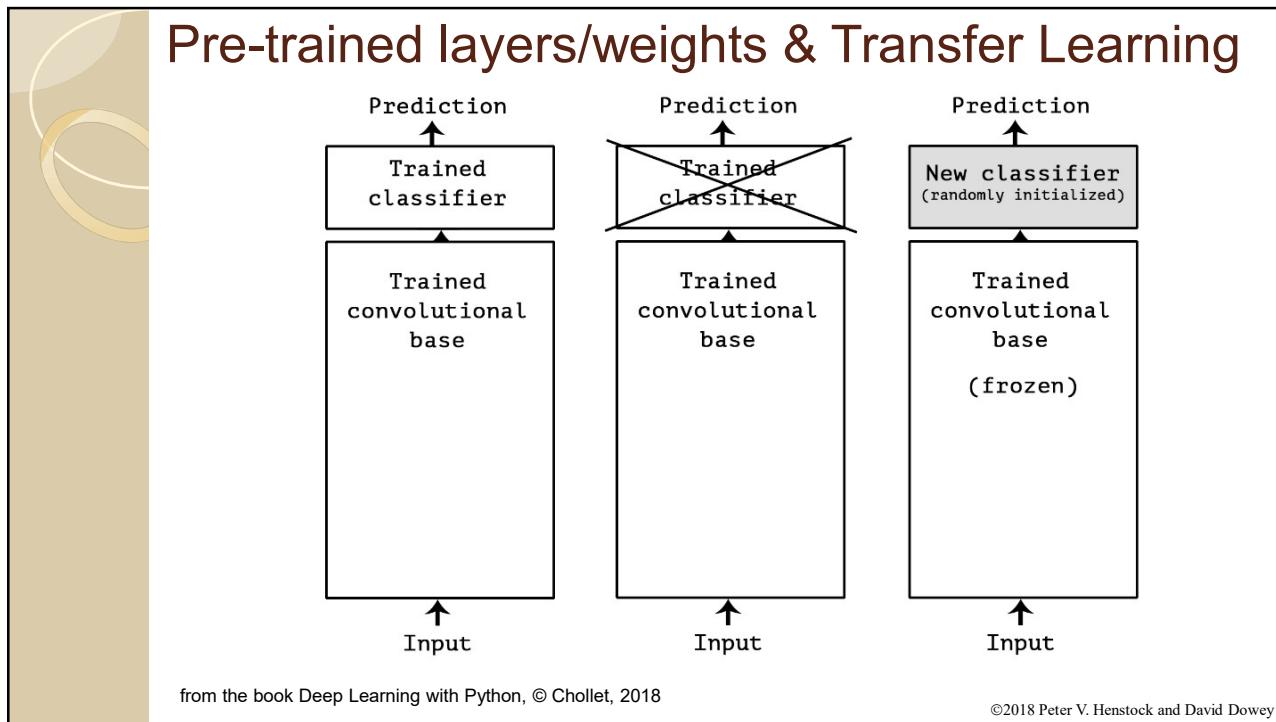
©2018 Peter V. Henstock and David Dowey

## Pre-trained layers/weights & Transfer Learning



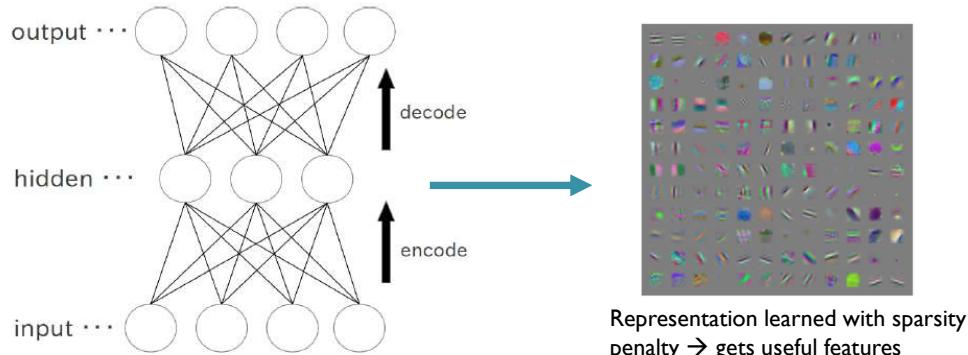
from the book Hands on Machine Learning with Scikit-Learn and TensorFlow, © Géron, 2017

©2018 Peter V. Henstock and David Dowey



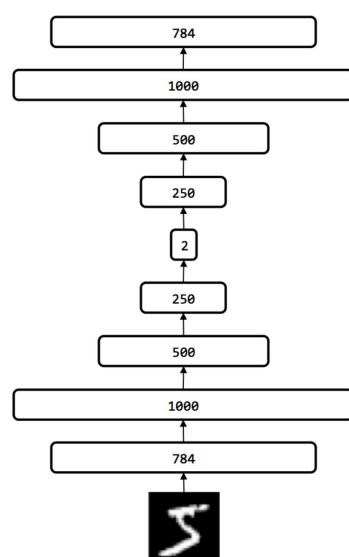
## Autoencoder Neural Network

- Way of compressing data; denoising
- Set input = output and learn (not identity function)
- Constrain # hidden nodes:  $\text{ceil}(\log_2(n))$  for n inputs
- Hidden layer learns the most salient pattern



©2018 Peter V. Henstock and David Dowey

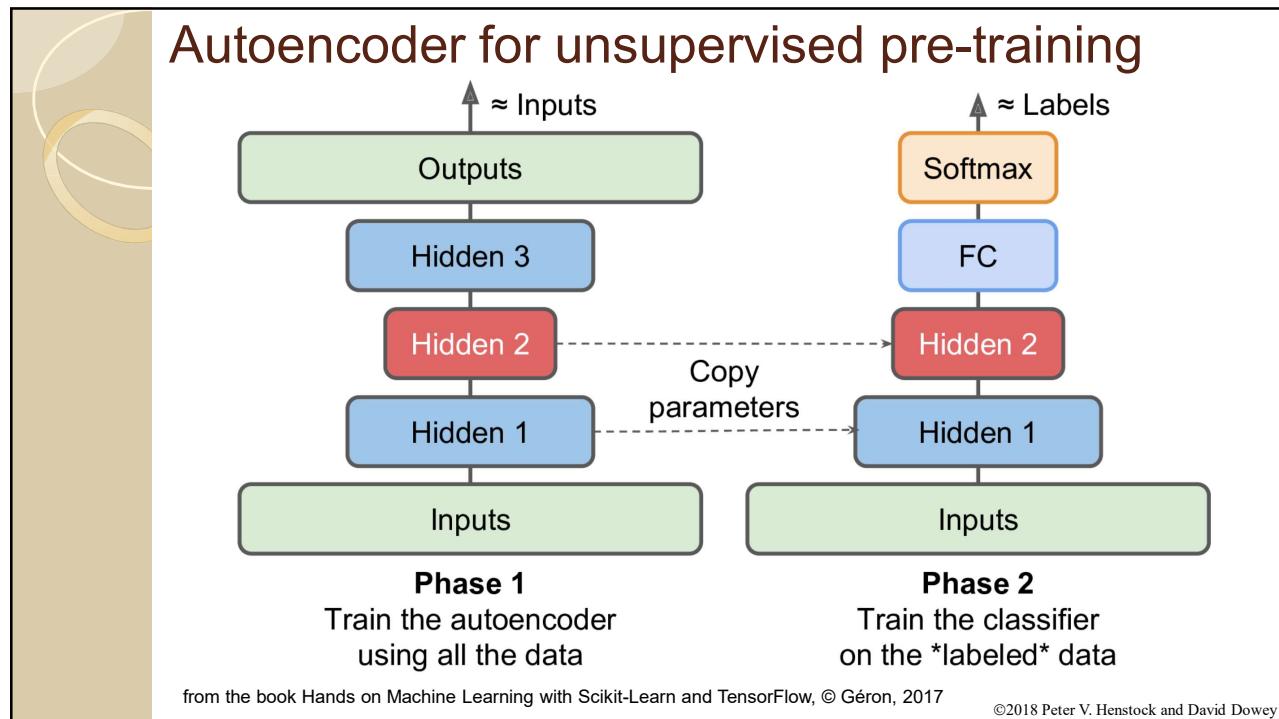
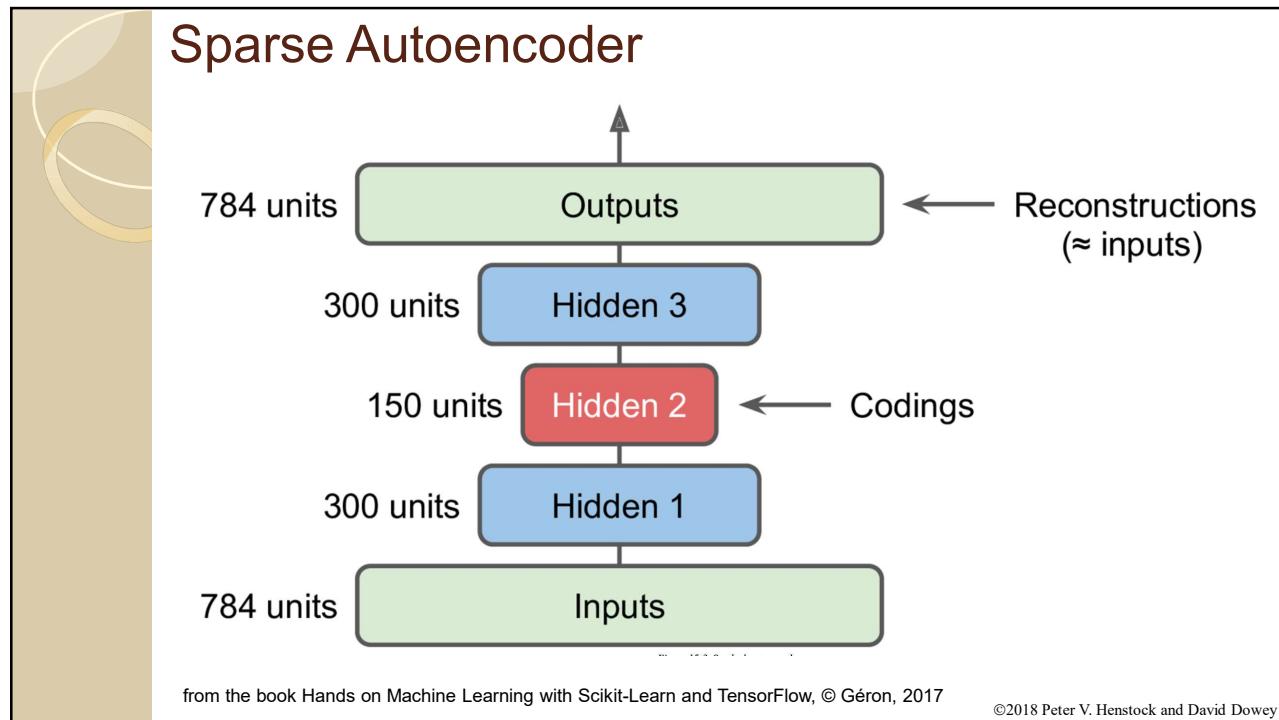
## Sparse Autoencoder



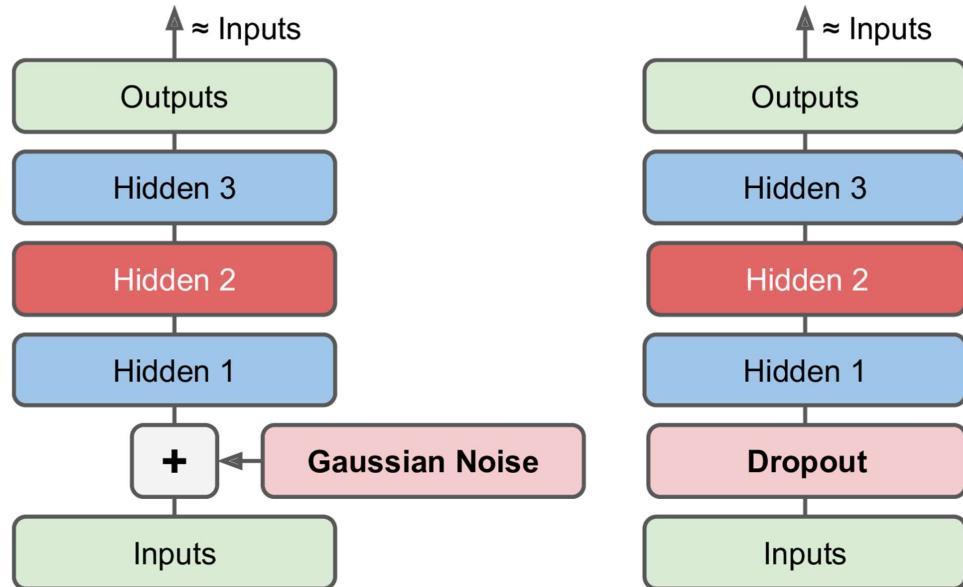
- Same image as input and output
- Learning to create the optimal sparse representation
- Here: MNIST digit

from the book Fundamentals of Deep Learning, © Buduma, 2017

©2018 Peter V. Henstock and David Dowey



## Autoencoder for denoising



from the book Hands on Machine Learning with Scikit-Learn and TensorFlow, © Géron, 2017

©2018 Peter V. Henstock and David Dowey