# Deep Learning 1

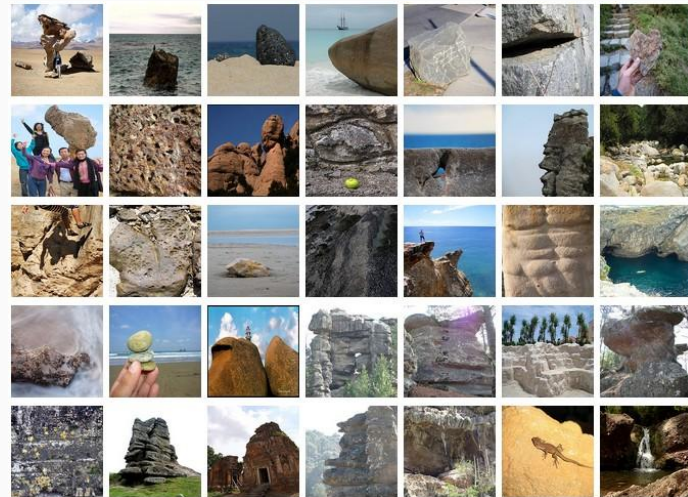# Pointers from Last Time

- Saving your models
    - Just need to save parameters

    - Can Use it anywhere

    - E.g. in linear regression with y=w1x1 + w2 you can just save w1 and w2 forever!

    - No data dependence after training

- Cost function is only used during training. Prediction function is eternal.

- Linear regression can be made non linear just by adding multiple powers of the same feature

# Motivation: ImageNet

# Image Classification Dataset

~1,000,000 images
~1,000 classes


Ground truths prepared manually
through Amazon Mechanical Turk

ImageNet Top-5 challenge:

You score if ground truth class is one your top 5 predictions

Best approaches used hand-crafted  features
(SIFT, HOGs, Fisher vectors, etc)
+ classifier

Top-5 error rate: ~25%

The Game Has Changed...

Krizhevsky, Sutskever and Hinton;
*ImageNet Classification with
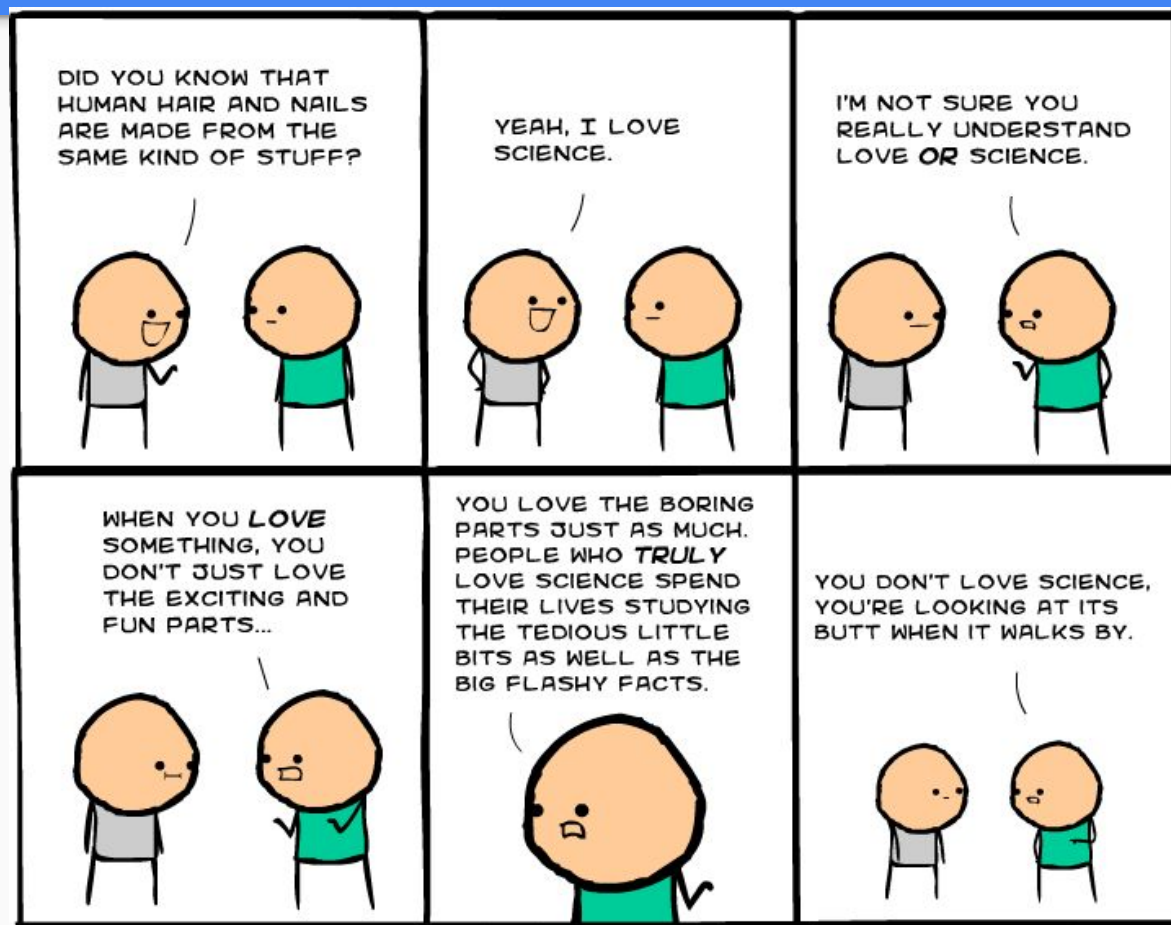Deep Convolutional Neural networks* **[Krizhevsky12]**

Top-5 error rate of ~15%

In the last few years, more modern networks have achieved better results.
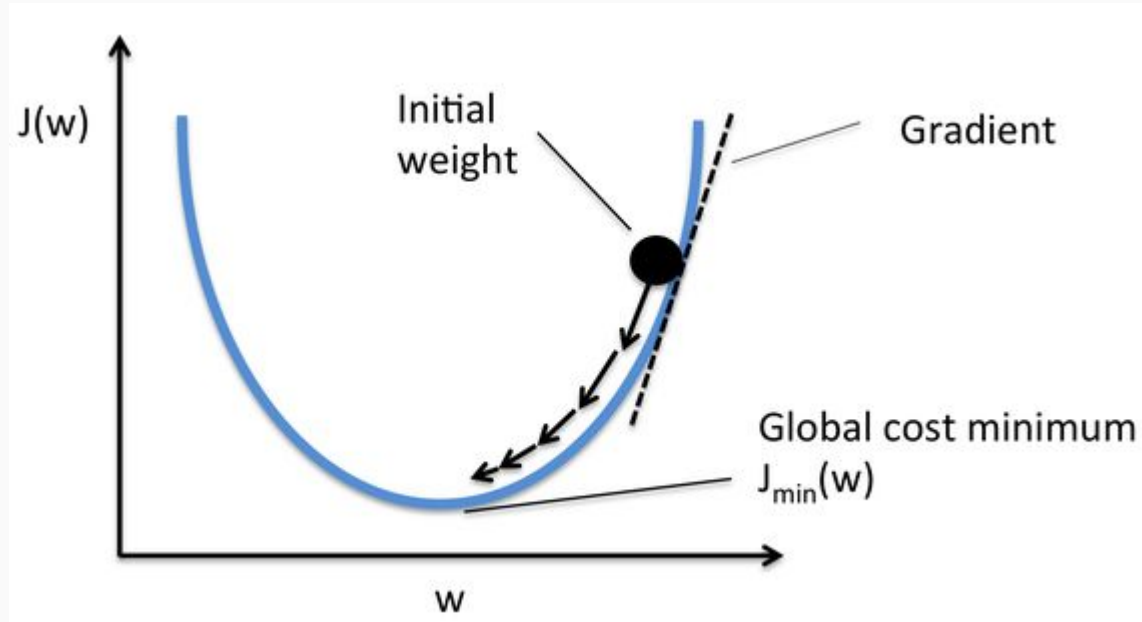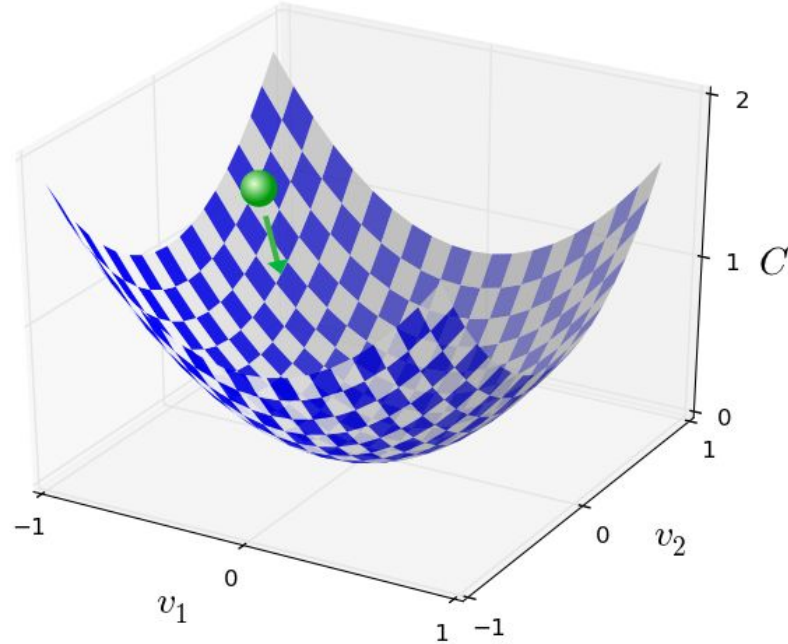**[Simonyan14, He15]**

Top-5 error rates of ~5-7%

- We talked about loss functions (i.e. cost functions)
  - Absolute Loss $| h(x) - y |$
  - Quadrative Loss $(h(x) - y)^2$

- Real world is not so generous.

- The following cycle goes on and on until we have a reasonable model
  - Predict using current hypothesis (last time)
  - Find how good the prediction was (last time)
  - **Update the hypothesis. (this time)**

- Q: How Does the update step happen?
  - A: Gradient Descent

# Gradient Descent (Cont.)

http://neuralnetworksanddeeplearning.com/images/valley_with_ball.png

# Algorithm

## Gradient descent algorithm

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$
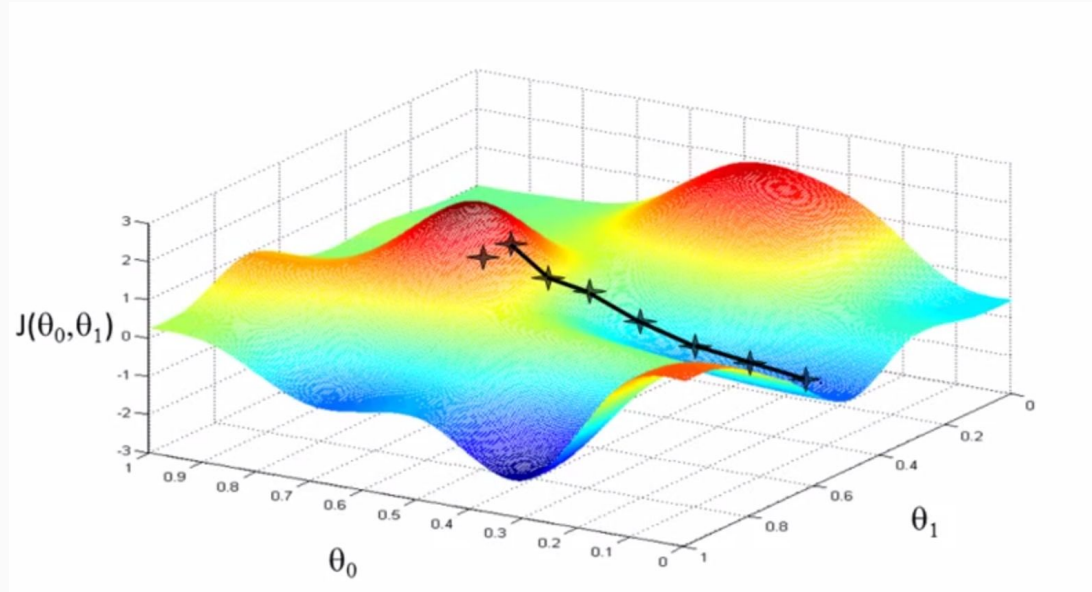
(for $j = 1$ and $j = 0$)

}

## Repeat until convergence

{

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

}

# Real world is ugly
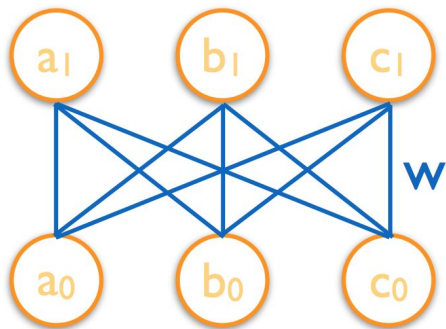


https://kousikk.files.wordpress.com/2014/11/screen-shot-2014-11-12-at-2-58-28-pm.png

TensorFlow

import tensorflow as tf

x            w

$$a_0 \quad b_0 \quad c_0 \; \cdot \; \begin{array}{|c|c|c|} \hline w_{a,a} & w_{a,b} & w_{a,c} \\ \hline w_{b,a} & w_{b,b} & w_{b,c} \\ \hline w_{c,a} & w_{c,b} & w_{c,c} \\ \hline \end{array} \; = \; a_1 \quad b_1 \quad c_1$$

y = tf.matmul(x, w)

$a_1$ =relu( $a_1$ )

$b_1$ =relu( $b_1$ )

$c_1$ =relu( $c_1$ )

out = tf.nn.relu(y)

# Data Flow Graph

- Describe mathematical computation with a directed graph of <u>nodes</u> & <u>edges</u>.
  - Nodes in the graph represent mathematical operations.

  - Edges describe the i/o relationships between nodes.

  - Data edges carry dynamically-sized multidimensional data arrays, or **tensors**.

- Nodes are assigned to computational devices and execute asynchronously and in parallel once all the tensors on their incoming edges becomes available.
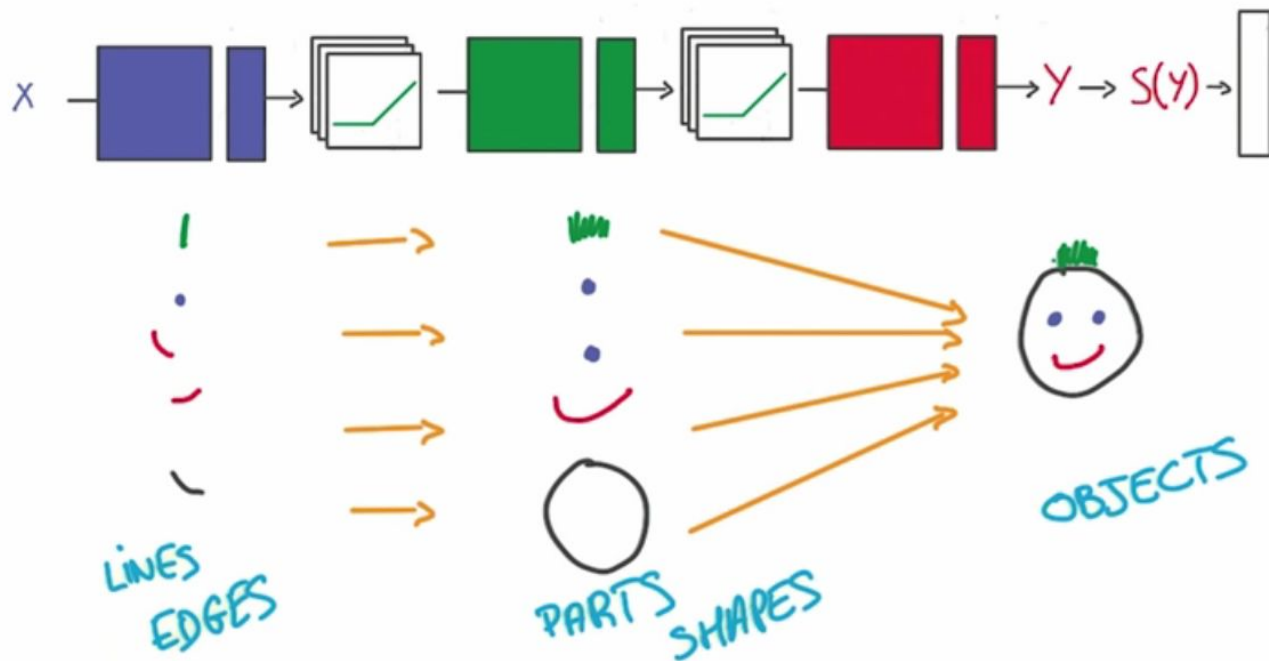
# TensorFlow Playground

http://playground.tensorflow.org/

# Back to Code: Hello Tensor

- ❏ Multiple layers

- ❏ Data propagates through layers

- ❏ Transformed by each layer

DEEP NETWORKS

$X \rightarrow$ ... $\rightarrow Y \rightarrow S(Y) \rightarrow$

LINES EDGES

PARTS SHAPES

OBJECTS

DeepLearning
udacity.com

# A Perceptron: (Logistic Regression)



http://www.parallelr.com/wp-content/uploads/2016/02/neuron.png

# Increasing The classifiers



http://pages.cs.wisc.edu/~bolo/shipyard/neural/nn_schematic.gif

# As a classifier



Input vector

Hidden layer 0 activation

Final layer activation (with softmax non-linearity)

Class probabilities

Image pixels

$P(cat) = 0.25$

$P(dog) = 0.5$

$P(car) = 0.1$

$P(ball) = 0.15$

...
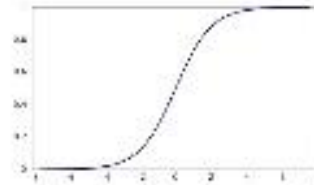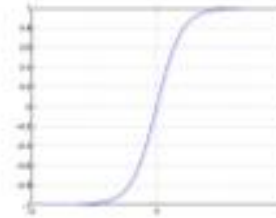
# Non-Linear Activation Function

- Sigmoid:   $S(t) = \dfrac{1}{1 + e^{-t}}$.

- Tanh:   $\tanh x = \dfrac{\sinh x}{\cosh x} = \dfrac{e^x - e^{-x}}{e^x + e^{-x}}$

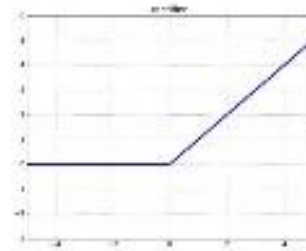- Rectified Linear Unit (ReLU):

  $f(x) = \max(0, x)$

Most popular activation function for DNN as of 2015, avoids saturation issues, makes learning faster
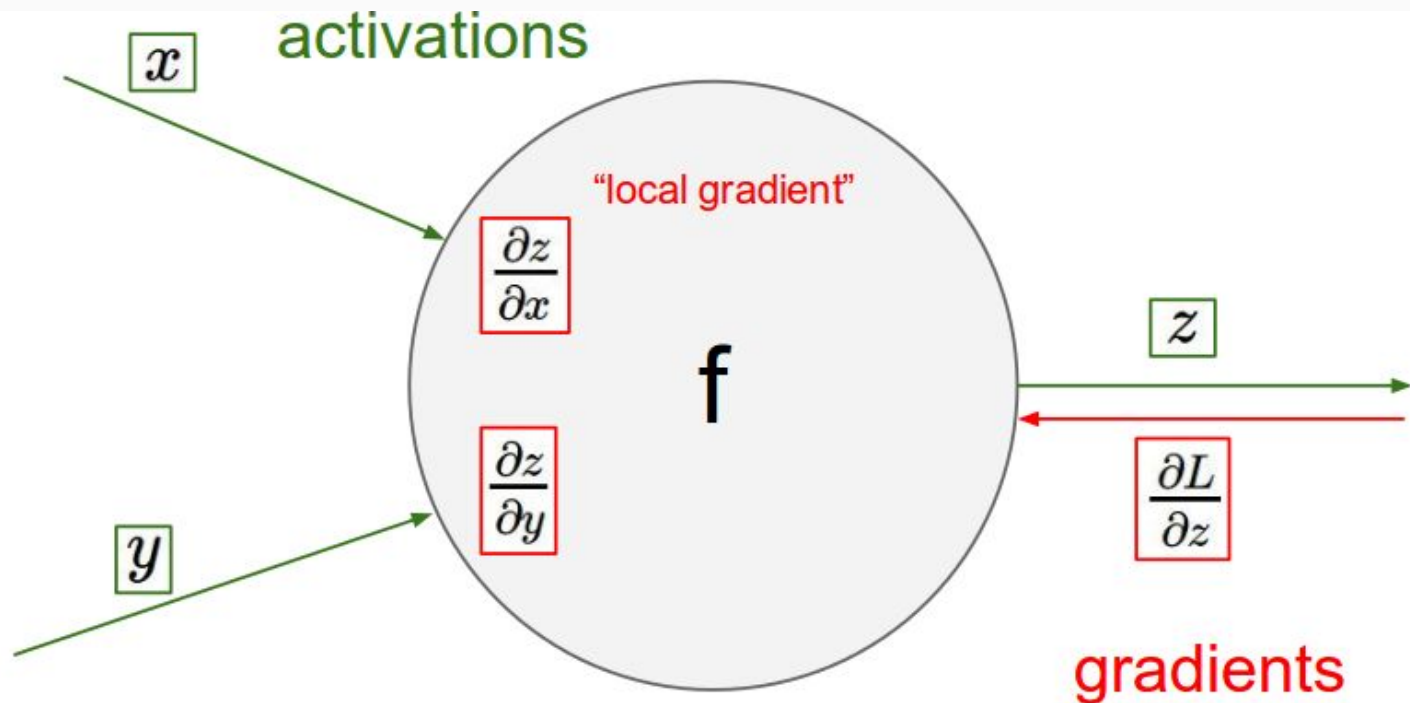
Sigmoid

Tanh

ReLU

# Training a Network

- Forward Propagation
  - Multiply Weights by inputs and add them up

  - Pass the result through a nonlinearity (ReLU, Sigmoid)

  - Repeat for next layer until the end layer

- Backward Propagation
  - Each node contributes to a certain degree of error

  - Final errors is a combination of everything

  - Disperse the error to respective nodes

  - Until you reach the starting node

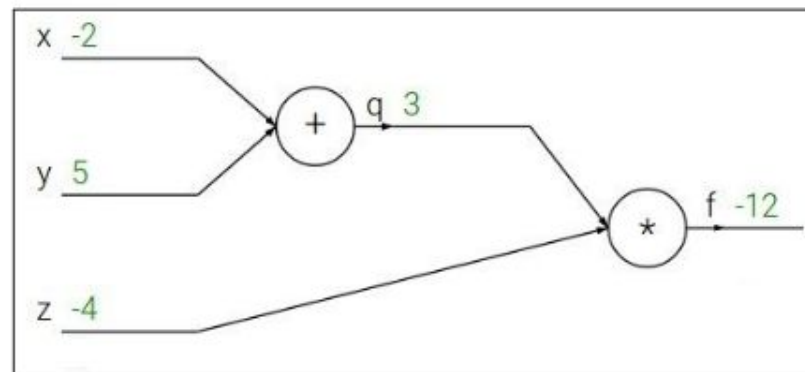$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

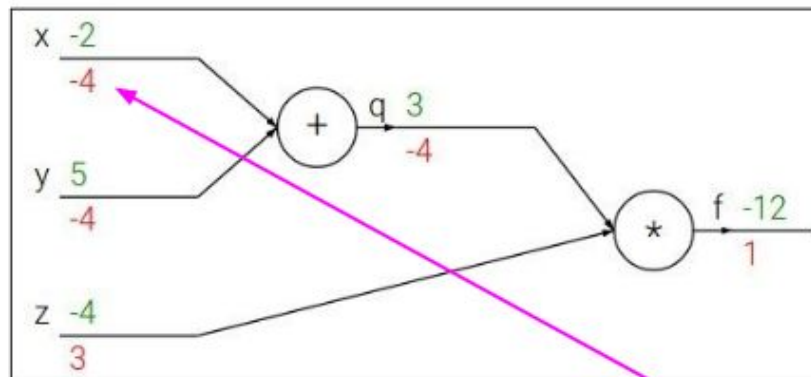Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$

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

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial x}$$

Chain rule:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$
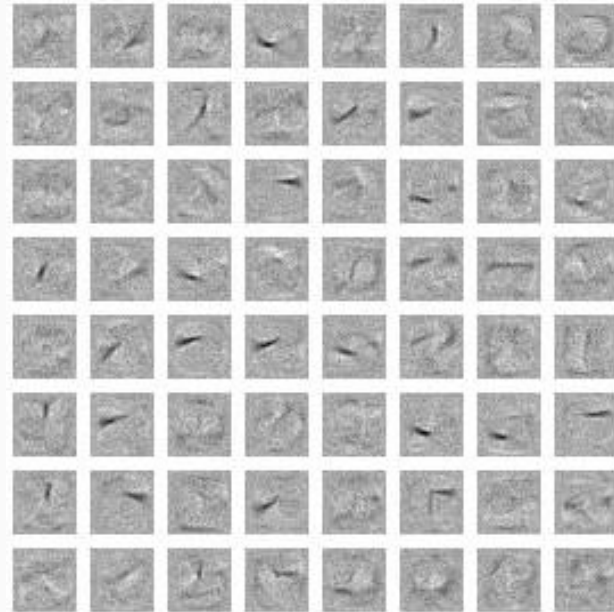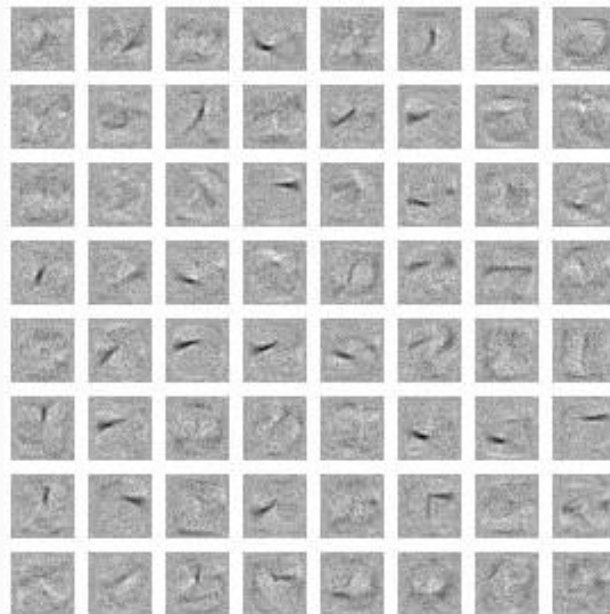
# Break (10 Min)

# Hands On After That

# A Problem

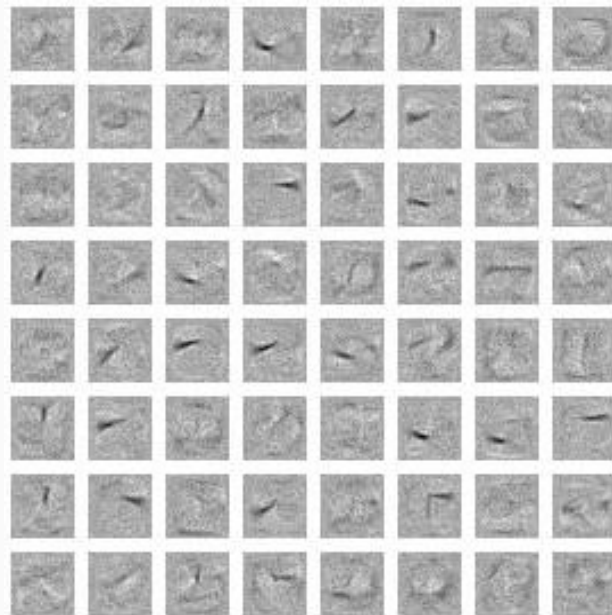# Visualising the learned weights can be educational

Each image visualises the weights connecting pixels
to a specific unit in the first hidden layer.

Note the stroke features detected by the various units

The fully connected networks so far have a weakness:

No translation invariance; learned features are
position dependent

For more general imagery:

- Requires a training set large enough to see all features in all possible positions…
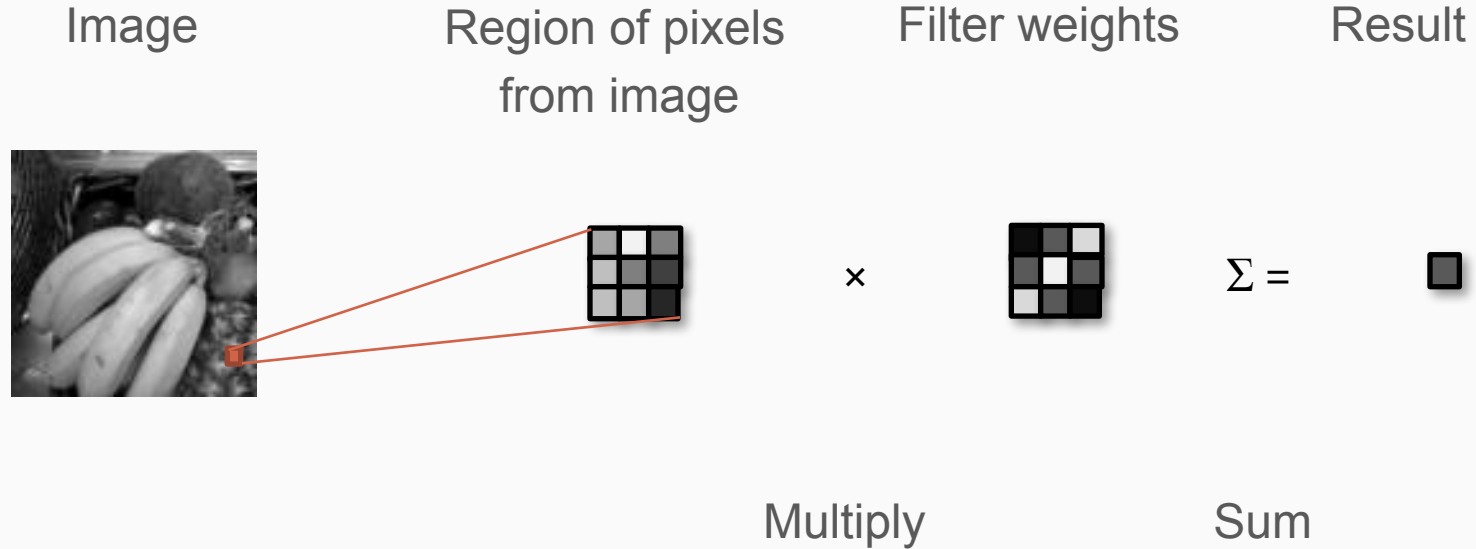
- Requires network with enough units to represent this…

# Convolutional Networks
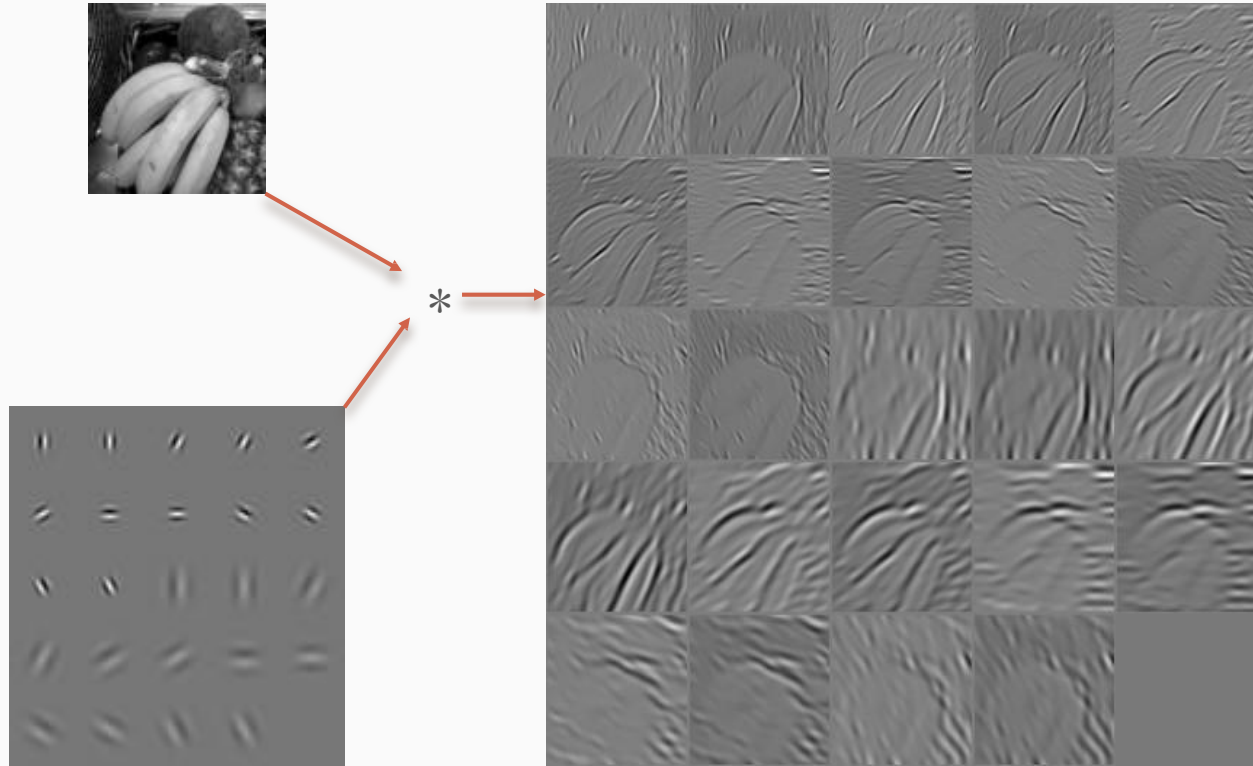
# Convolution
Often used for feature detection

Slide a convolutional filter over an image...

# Multiply image pixels by filter weights and sum

Image     Region of pixels from image     Filter weights     Result



×     Σ =

Multiply     Sum

## Do this for all possible positions in the image
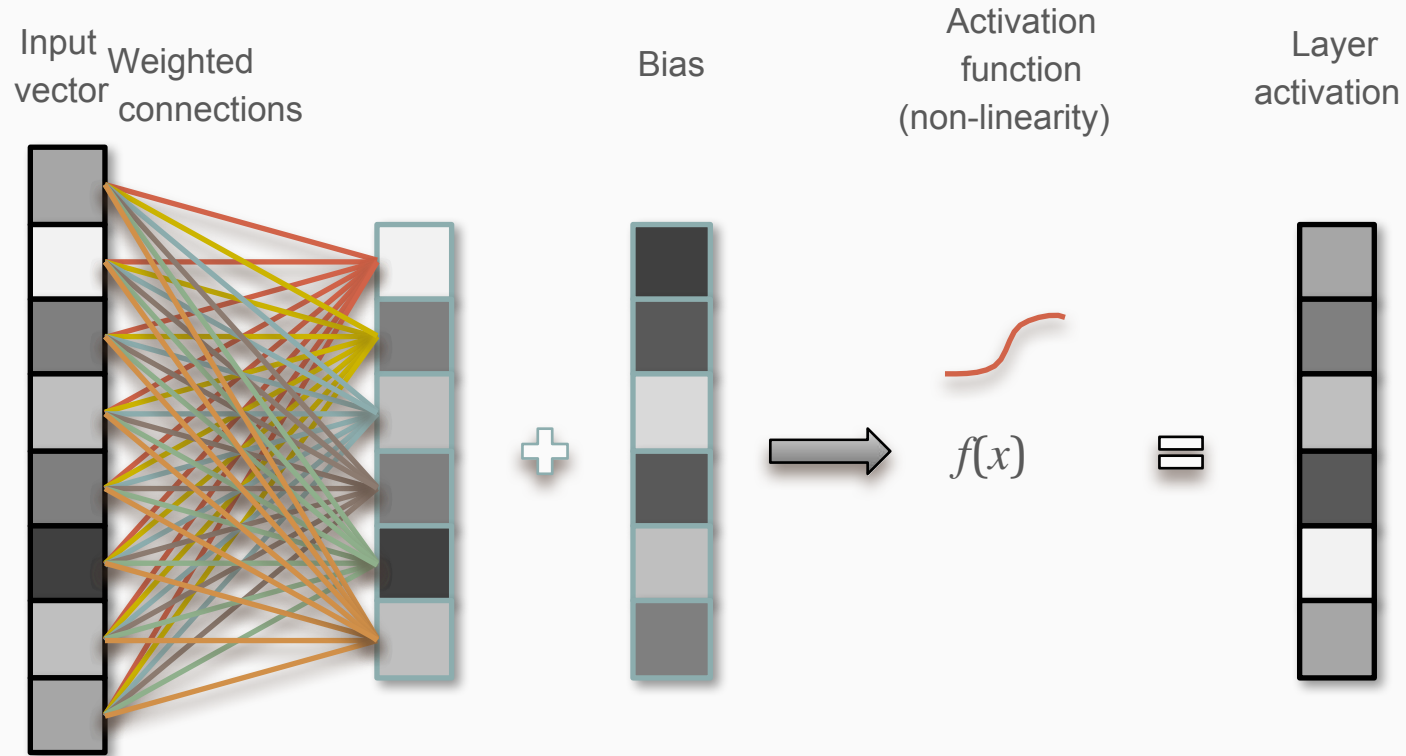
# Convolution: Gabor filters



*

Convolution detects features in a position independent manner
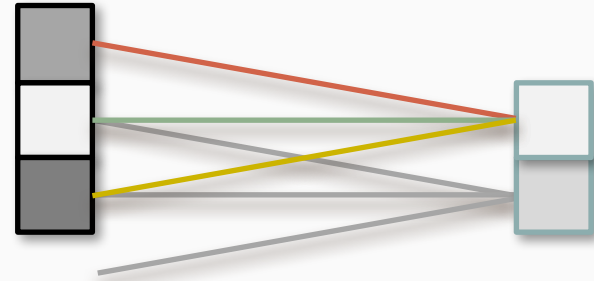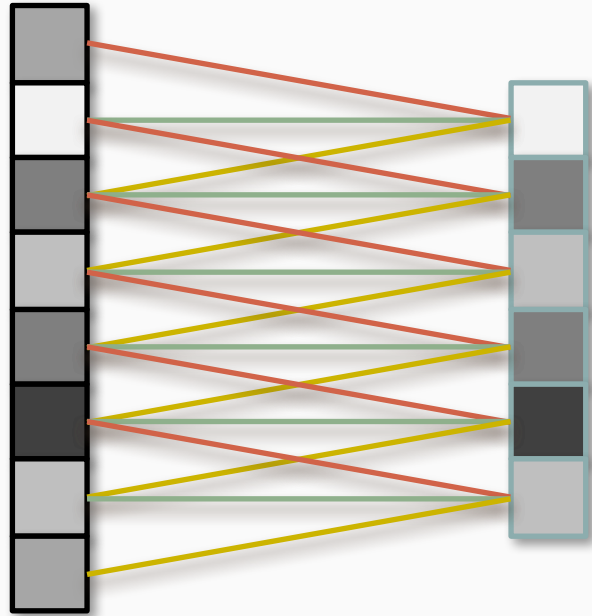
--

Convolutional neural networks learn position independent filters (feature detectors)

# Recap: FC (fully-connected) layer



Input vector  Weighted connections     Bias          Activation function (non-linearity)          Layer activation
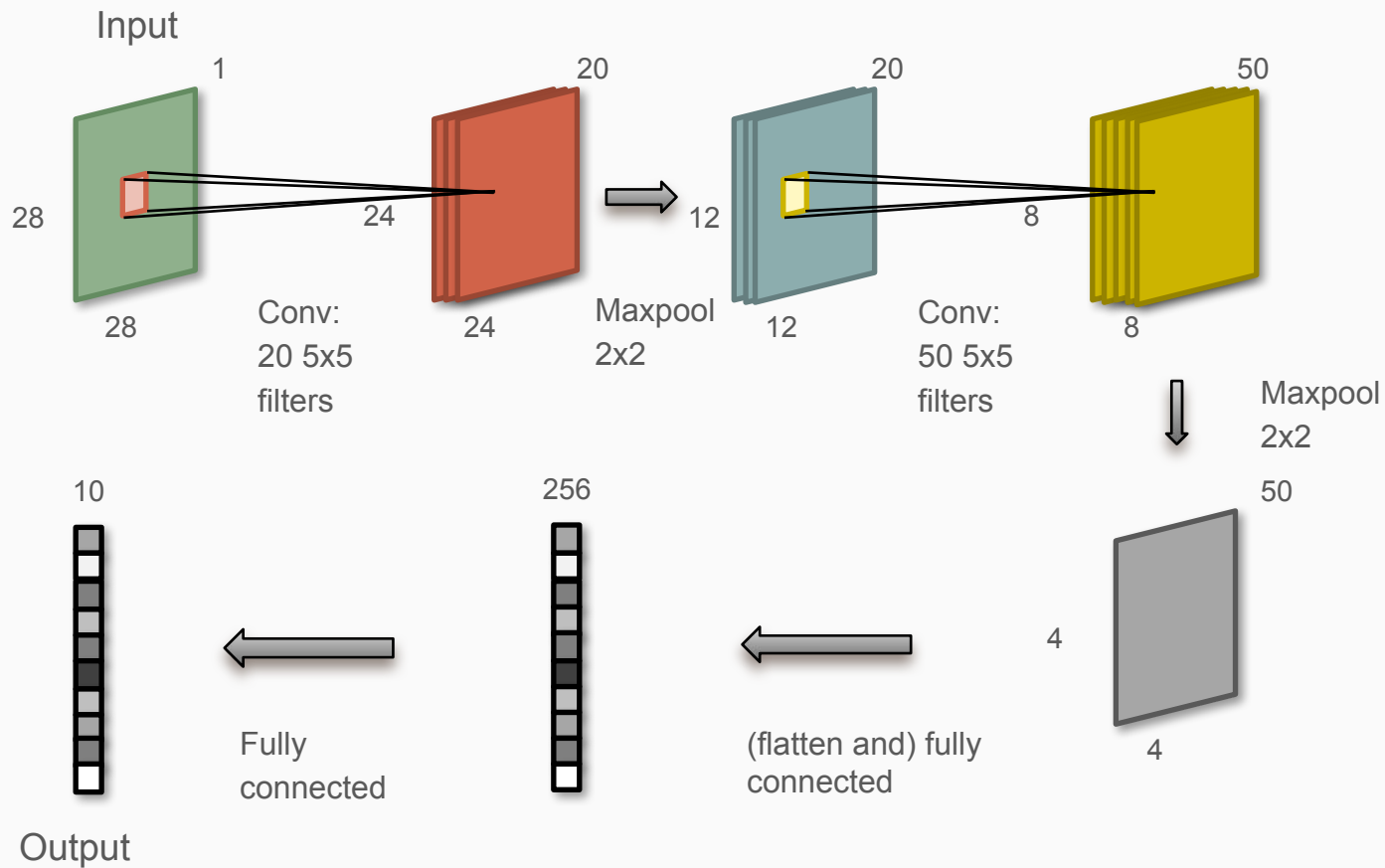
$f(x)$

# Convolutional layer



Each unit only connected to units in its neighbourhood

Example:


A Simplified LeNet **[LeCun95]** for MNIST digits

# Convolutional layer

Input

1

28

28

Conv:
20 5x5
filters

20

24

24

Maxpool
2x2

20

12

12

Conv:
50 5x5
filters

50

8

8

Maxpool
2x2

50

4

4

(flatten and) fully
connected

256

Fully
connected

10

Output

# Max-pooling 'layer' [Ciresan12]

Take maximum value from each 2 x 2  pooling region ($p$ x $p$) in the general case

after 300 iterations over training set:
**99.21%** validation accuracy

# Thanks

Adapted from:
- https://github.com/Britefury
- https://github.com/nlintz/TensorFlow-Tutorials