

Machine learning - HT 2016

8. Neural Networks

Varun Kanade

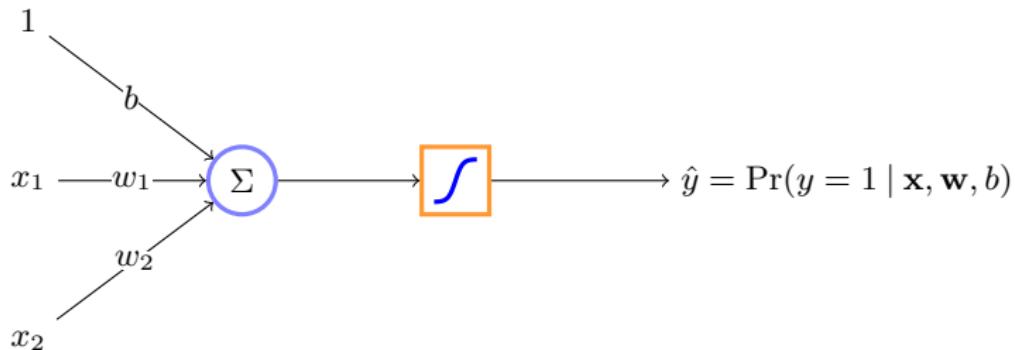
University of Oxford
February 19, 2016

Outline

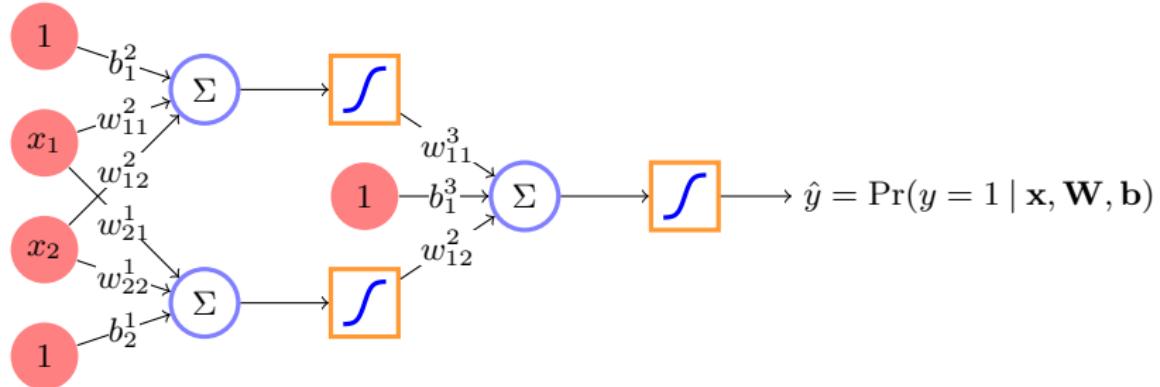
Today, we'll study feedforward neural networks

- ▶ Multi-layer perceptrons
- ▶ Application to classification or regression settings
- ▶ Backpropagation to compute gradients
- ▶ Finally understand what `model:forward` and `model:backward` is doing

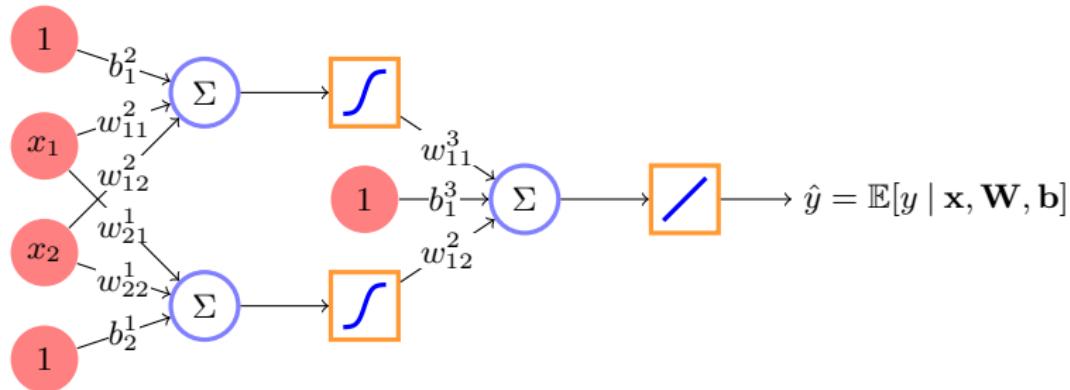
Artificial Neuron : Logistic Regression



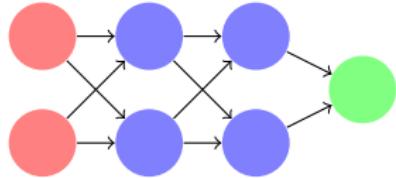
Multilayer Perceptron (MLP) : Classification



Multilayer Perceptron (MLP) : Regression



A Simple Example



$$\hat{y} = \sigma \left(w_{11}^4 \sigma \left(w_{11}^3 \sigma (w_{11}^2 x_1 + w_{12}^2 x_2 + b_1^2) + w_{12}^3 \sigma (w_{21}^2 x_1 + w_{22}^2 x_2 + b_2^2) + b_1^3 \right) \right. \\ \left. + w_{12}^4 \sigma \left(w_{21}^3 \sigma (w_{11}^2 x_1 + w_{12}^2 x_2 + b_1^2) + w_{22}^3 \sigma (w_{21}^2 x_1 + w_{22}^2 x_2 + b_2^2) + b_1^3 \right) + b_1^4 \right)$$

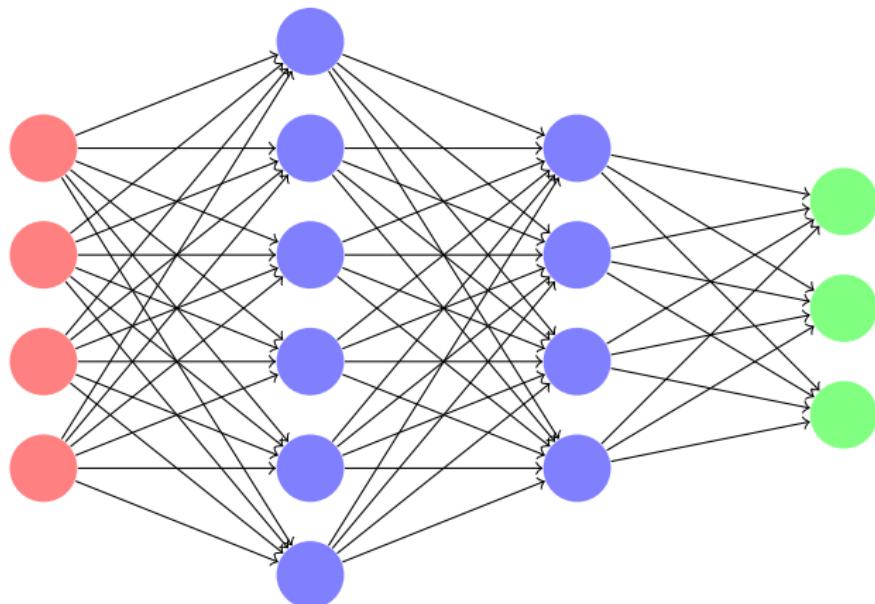
Feedforward Neural Networks

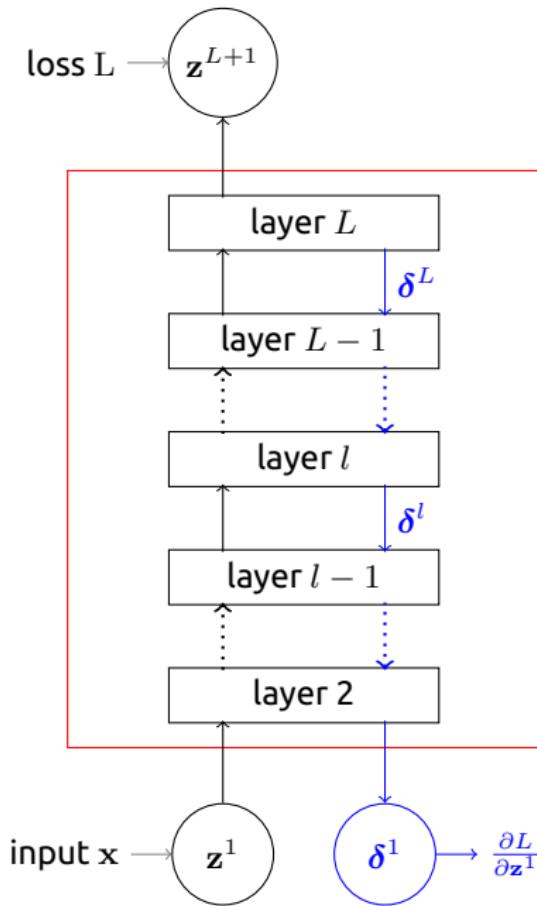
Layer 1
(Input)

Layer 2
(Hidden)

Layer 3
(Hidden)

Layer 4
(Output)



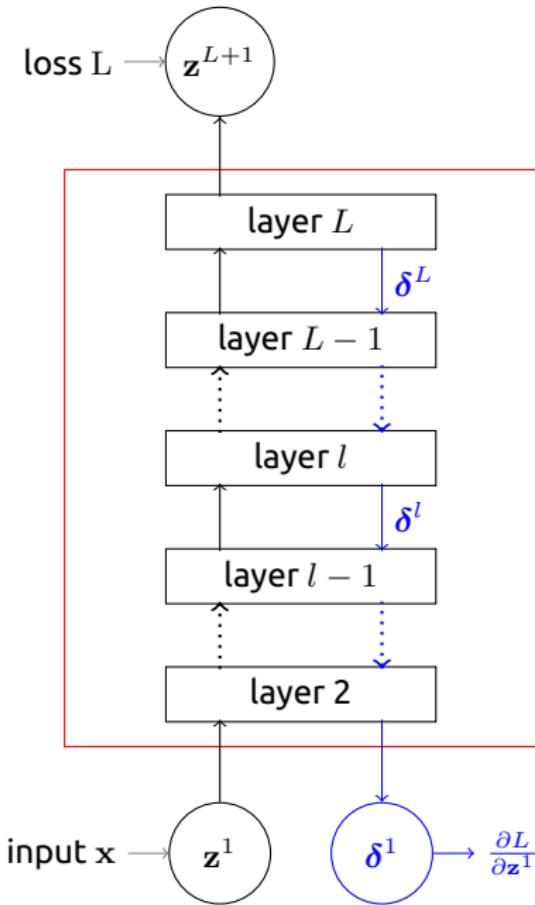


```

model = nn.Sequential()
model:add(nn.Linear(n_in, n_h1));
model:add(nn.Sigmoid());
...
...
model:add(nn.Linear(n_h_L-1, n_out));
model:Sigmoid();
criterion = nn.MSECriterion();

model:forward(x)
criterion:backward(model.output, y)

dL_do = criterion:backward(model.output,y)
model:backward(x, dL_do)
    
```



Forward Equations

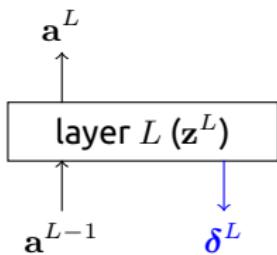
$$(1) \quad \mathbf{z}^1 = \mathbf{x} \text{ (input)}$$

$$(2) \quad \mathbf{z}^l = \mathbf{W}^l \mathbf{a}^{l-1} + \mathbf{b}^l$$

$$(3) \quad \mathbf{a}^l = f(\mathbf{z}^l)$$

$$(4) \quad L(\mathbf{a}^L, y)$$

Output Layer



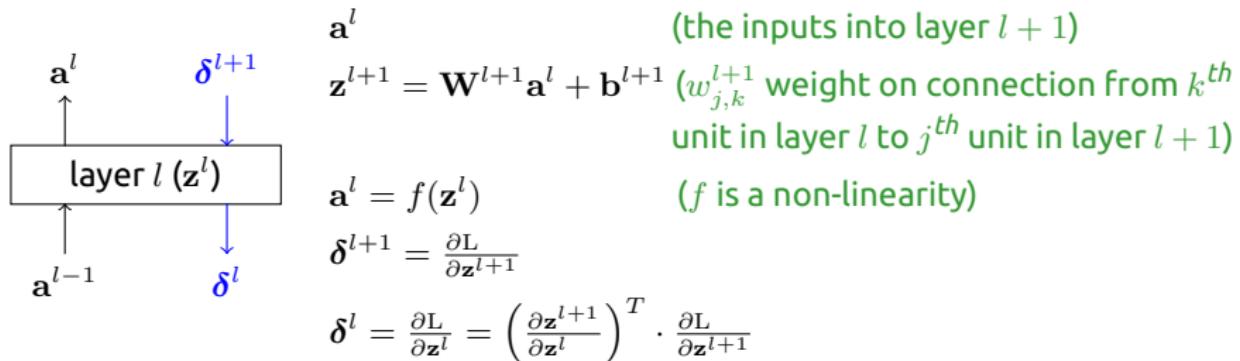
$$\mathbf{z}^L = \mathbf{W}^L \mathbf{a}^{L-1} + \mathbf{b}^L$$

$$\mathbf{a}^L = f(\mathbf{z}^L)$$

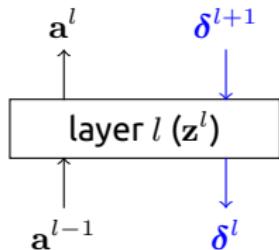
Loss: $L(y, \mathbf{a}^L) = \text{criterion : forward}(\mathbf{a}^L, y)$

$$\delta^L = \frac{\partial L}{\partial \mathbf{z}^L} = \frac{\partial L}{\partial \mathbf{a}^L} \odot f'(\mathbf{z}^L)$$

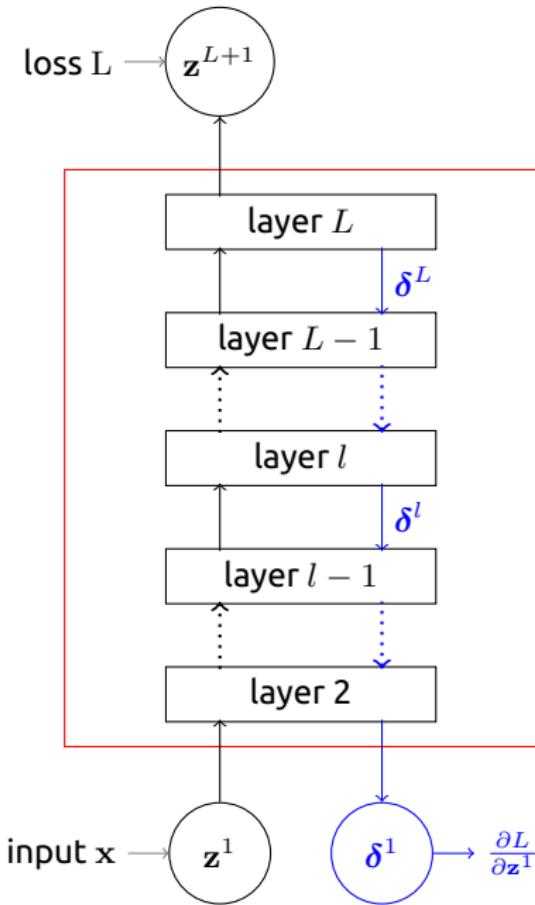
Back Propagation



Gradients wrt Parameters



$$\mathbf{z}^l = \mathbf{W}^l \mathbf{a}^{l-1} + \mathbf{b}^l \quad (\text{w}_{j,k}^l \text{ weight on connection from } k^{\text{th}} \text{ unit in layer } l-1 \text{ to } j^{\text{th}} \text{ unit in layer } l)$$
$$\delta^l = \frac{\partial L}{\partial \mathbf{z}^l}$$



Forward Equations

$$(1) \quad \mathbf{z}^1 = \mathbf{x} \text{ (input)}$$

$$(2) \quad \mathbf{z}^l = \mathbf{W}^l \mathbf{a}^{l-1} + \mathbf{b}^l$$

$$(3) \quad \mathbf{a}^l = f(\mathbf{z}^l)$$

$$(4) \quad L(\mathbf{a}^L, y)$$

Back-propagation Equations

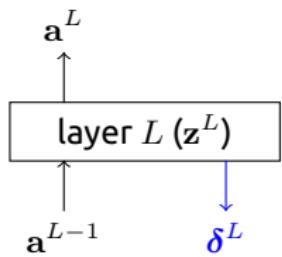
$$(1) \quad \delta^L = \frac{\partial L}{\partial \mathbf{a}^L} \odot f'(\mathbf{z}^L)$$

$$(2) \quad \delta^l = (\mathbf{W}^{l+1})^T \delta^{l+1} \odot f'(\mathbf{z}^l)$$

$$(3) \quad \frac{\partial L}{\partial \mathbf{b}^l} = \delta^l$$

$$(4) \quad \frac{\partial L}{\partial \mathbf{W}^l} = \delta^l (\mathbf{a}^{l-1})^T$$

Output Layer : LogSoftmax

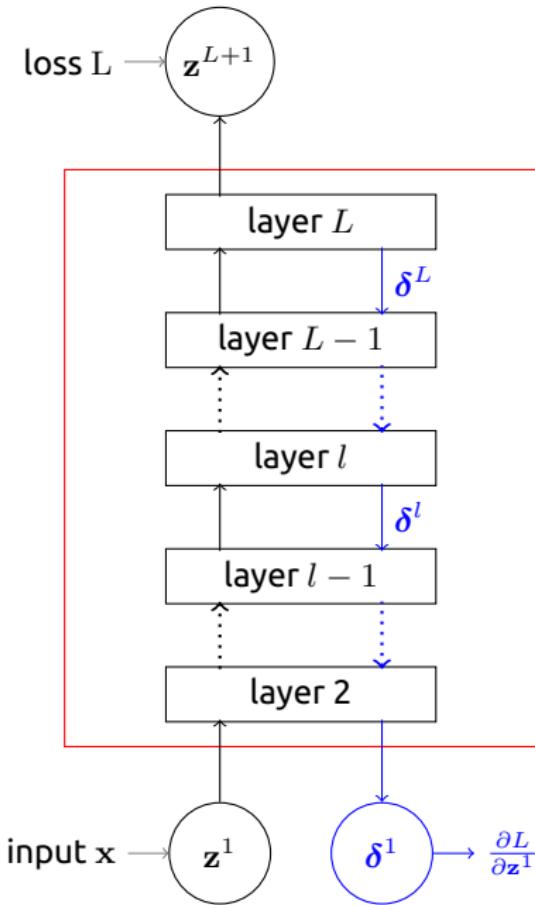


```
criterion = nn.ClassNLLCriterion();
```

Loss: $L(y, \mathbf{a}^L) = -a_y^L$

$$\mathbf{a}^L = \text{LogSoftmax}(\mathbf{z}^L)$$

$$\delta^L = \frac{\partial L}{\partial \mathbf{z}^L} = \left(\frac{\partial \mathbf{a}^L}{\partial \mathbf{z}^L} \right)^T \cdot \frac{\partial L}{\partial \mathbf{a}^L}$$



Forward Equations

$$(1) \quad \mathbf{z}^1 = \mathbf{x} \text{ (input)}$$

$$(2) \quad \mathbf{z}^l = g^l(\mathbf{a}^{l-1}; \mathbf{W}^l, \mathbf{b}^l)$$

$$(3) \quad \mathbf{a}^l = f^l(\mathbf{z}^l)$$

$$(4) \quad L(\mathbf{a}^L, y)$$

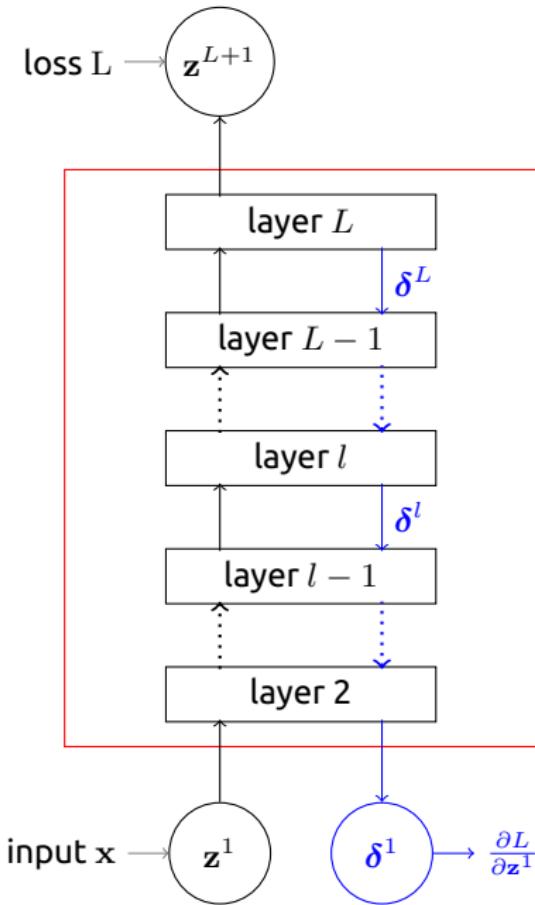
Back-propagation Equations

$$(1) \quad \delta^L = \left(\frac{\partial \mathbf{a}^L}{\partial \mathbf{z}^L} \right)^T \cdot \frac{\partial L}{\partial \mathbf{a}^L}$$

$$(2) \quad \delta^l = \left(\frac{\partial \mathbf{z}^{l+1}}{\partial \mathbf{a}^l} \cdot \frac{\partial \mathbf{a}^l}{\partial \mathbf{z}^l} \right)^T \cdot \delta^{l+1}$$

$$(3) \quad \frac{\partial L}{\partial \mathbf{b}^l} = \left(\frac{\partial \mathbf{z}^l}{\partial \mathbf{b}^l} \right)^T \cdot \delta^l$$

$$(4) \quad \frac{\partial L}{\partial \mathbf{W}^l} = \left(\frac{\partial \mathbf{z}^l}{\partial \mathbf{W}^l} \right)^T \cdot \delta^l$$



Forward Equations

$$(1) \quad \mathbf{z}^1 = \mathbf{x} \text{ (input)}$$

$$(2) \quad \mathbf{z}^l = \mathbf{W}^l \mathbf{a}^{l-1} + \mathbf{b}^l$$

$$(3) \quad \mathbf{a}^l = f(\mathbf{z}^l)$$

$$(4) \quad L(\mathbf{a}^L, y)$$

Back-propagation Equations

$$(1) \quad \delta^L = \frac{\partial L}{\partial \mathbf{a}^L} \odot f'(\mathbf{z}^L)$$

$$(2) \quad \delta^l = (\mathbf{W}^{l+1})^T \delta^{l+1} \odot f'(\mathbf{z}^l)$$

$$(3) \quad \frac{\partial L}{\partial \mathbf{b}^l} = \delta^l$$

$$(4) \quad \frac{\partial L}{\partial \mathbf{W}^l} = \delta^l (\mathbf{a}^{l-1})^T$$

Computational Questions

What is the running time to compute the gradient for a single data point?

What is the space requirement?

Can we process multiple examples together?

```
model = nn.Sequential()
model.add(nn.Linear(3, 4))
model.add(nn.Sigmoid())
model.add(nn.Linear(4, 2))
model.add(nn.Sigmoid())
criterion = nn.MSELoss()
```

```
: model.modules
: {
  1 :
    nn.Linear(3 -> 4)
    {
      gradBias : DoubleTensor - size: 4
      weight : DoubleTensor - size: 4x3
      gradWeight : DoubleTensor - size: 4x3
      gradInput : DoubleTensor - empty
      bias : DoubleTensor - size: 4
      output : DoubleTensor - empty
    }
  2 :
    nn.Sigmoid
    {
      gradInput : DoubleTensor - empty
      output : DoubleTensor - empty
    }
  3 :
    nn.Linear(4 -> 2)
    {
      gradBias : DoubleTensor - size: 2
      weight : DoubleTensor - size: 2x4
      gradWeight : DoubleTensor - size: 2x4
      gradInput : DoubleTensor - empty
      bias : DoubleTensor - size: 2
      output : DoubleTensor - empty
    }
  4 :
    nn.Sigmoid
    {
      gradInput : DoubleTensor - empty
      output : DoubleTensor - empty
    }
}
```

```
model = nn.Sequential()
model.add(nn.Linear(3, 4))
model.add(nn.Sigmoid())
model.add(nn.Linear(4, 2))
model.add(nn.Sigmoid())
criterion = nn.MSELoss()
inputs = torch.randn(500, 3)
targets = torch.randn(500, 2)
model.forward(inputs)
criterion.forward(model.output, targets)
deriv = criterion.backward(model.output, targets)
model.backward(inputs, deriv)
```

```
model.modules
{
    1 :
        nn.Linear(3 -> 4)
    {
        gradBias : DoubleTensor - size: 4
        weight : DoubleTensor - size: 4x3
        gradWeight : DoubleTensor - size: 4x3
        gradInput : DoubleTensor - size: 500x3
        addBuffer : DoubleTensor - size: 500
        bias : DoubleTensor - size: 4
        output : DoubleTensor - size: 500x4
    }
    2 :
        nn.Sigmoid
    {
        gradInput : DoubleTensor - size: 500x4
        output : DoubleTensor - size: 500x4
    }
    3 :
        nn.Linear(4 -> 2)
    {
        gradBias : DoubleTensor - size: 2
        weight : DoubleTensor - size: 2x4
        gradWeight : DoubleTensor - size: 2x4
        gradInput : DoubleTensor - size: 500x4
        addBuffer : DoubleTensor - size: 500
        bias : DoubleTensor - size: 2
        output : DoubleTensor - size: 500x2
    }
    4 :
        nn.Sigmoid
    {
        gradInput : DoubleTensor - size: 500x2
        output : DoubleTensor - size: 500x2
    }
}
```

Training Deep Neural Networks

- ▶ Back-propagation gives gradient
- ▶ Stochastic gradient descent is the method of choice
- ▶ Regularisation
 - ▶ How do we add ℓ_1 or ℓ_2 regularisation?
 - ▶ Don't regularise bias terms
- ▶ How about convergence?
- ▶ What did we learn in the last 10 years, that we didn't know in the 80s?

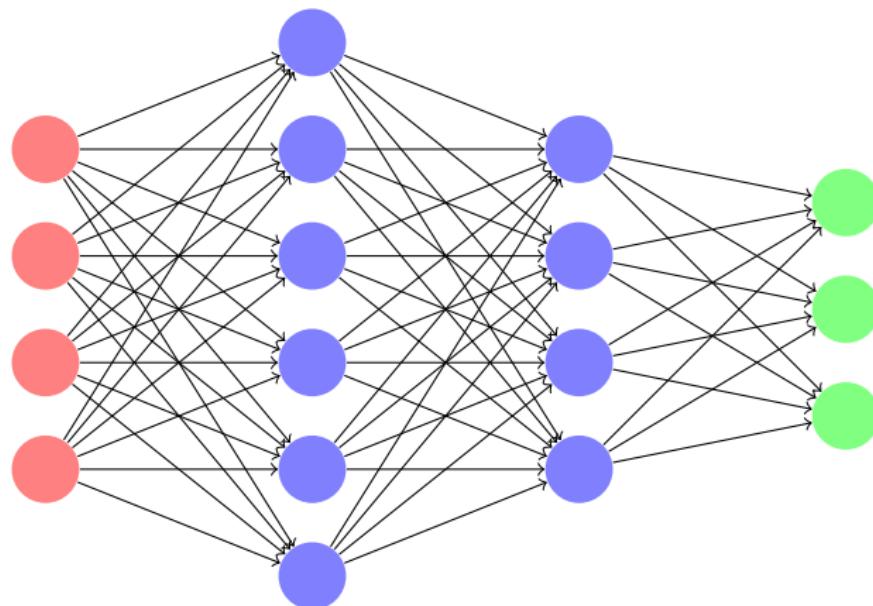
Training Feedforward Deep Networks

Layer 1
(Input)

Layer 2
(Hidden)

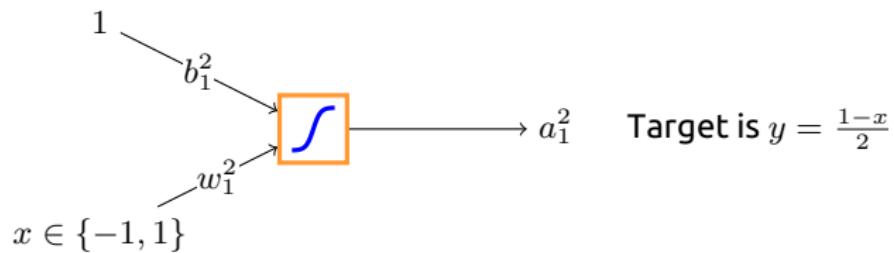
Layer 3
(Hidden)

Layer 4
(Output)

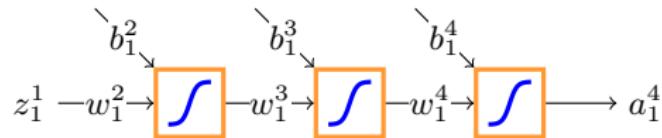


Why do we get non-convex optimisation problem?

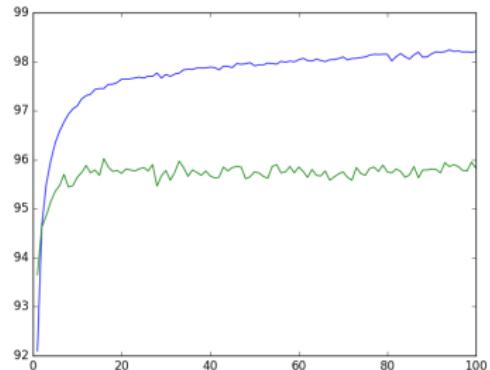
A toy example



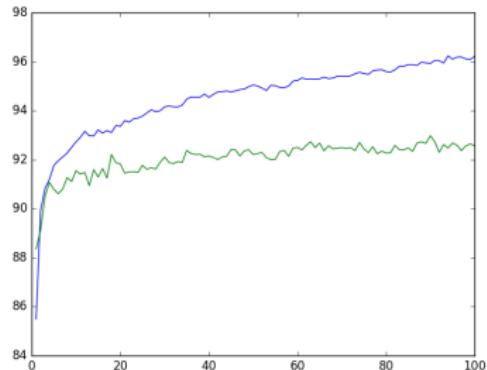
Propagating Gradients Backwards



Digit Classification Problem on Sheet 4

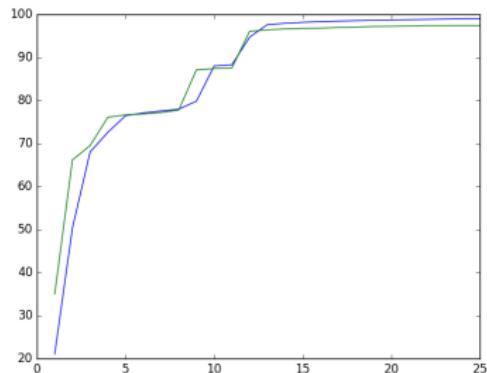


One-hot encoding

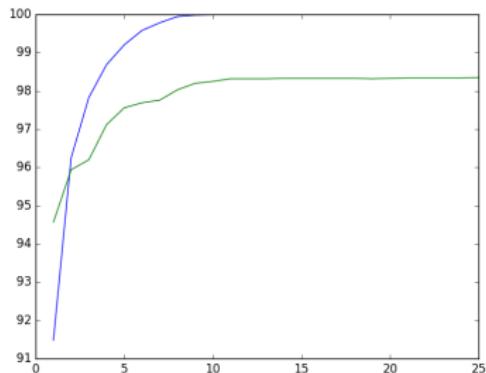


Binary encoding

Digit Classification: Squared Loss vs Cross Entropy



Squared Error



Cross Entropy

(See paper by Glorot and Bengio (2010))

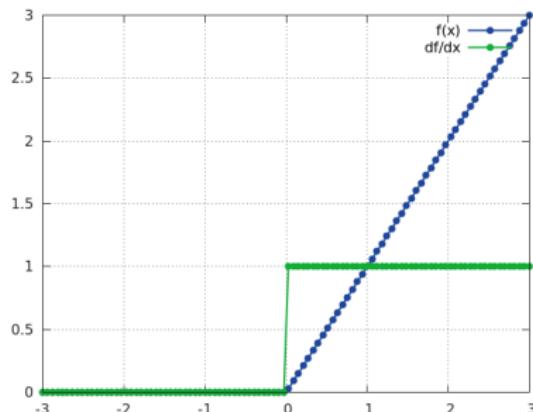
Avoiding Saturation

Use *rectified linear units*

$$\text{ReLU}(z) = \max(0, z)$$

Sometimes, you will see just
 $f(z) = |z|$

Other varieties: leaky ReLUs,
parametric ReLUs



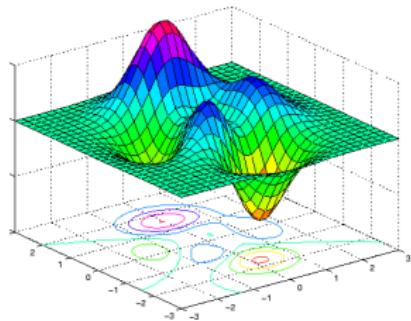
Initialising Weights and Biases

Why is initialising important?

Suppose we were using a *sigmoid unit*, how would you initialise the incoming weights?

What if it were a ReLU unit?

How about the biases?



Avoiding Overfitting

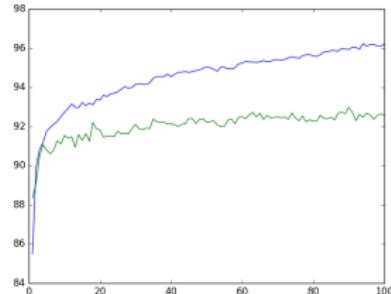
Deep Neural Networks have a lot of parameters

- ▶ Fully connected layers with n_1, n_2, \dots, n_k units have at least $n_1 n_2 + n_2 n_3 + \dots + n_{k-1} n_k$ parameters
- ▶ MLP for digit recognition had 2 million parameters!
- ▶ For image detection, the neural net used by Krizhevsky, Sutskever, Hinton (2012) has 60 *million* parameters and 1.2 *million* training images
- ▶ How do we avoid deep neural networks from overfitting?

Early Stopping

Maintain validation set and stop training when error on validation set stops decreasing.

What are the computational costs?



What are the advantages?

Early stopping leads to better generalisation

(See paper by Hardt, Recht and Singer (2015))

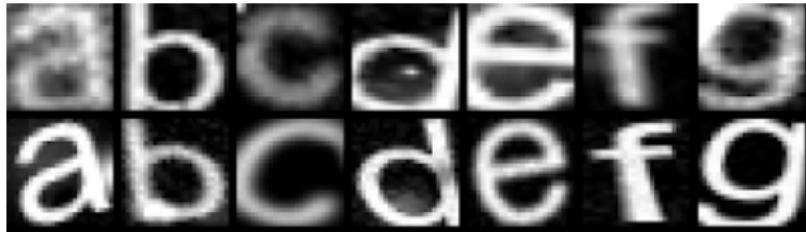
Add Data

Typically, getting additional data is either impossible or expensive

Fake the data!

Images can be translated slight, rotated slightly, change of brightness, etc.

Google Offline Translate trained on entirely fake data!



(Google Research Blog)

Add Data

Adversarial Training

Take trained (or partially trained model)

Create examples by modifications “imperceptible to the human eye”, but where the model fails

$$\mathbf{x} + .007 \times \text{sign}(\nabla_{\mathbf{x}} J(\boldsymbol{\theta}, \mathbf{x}, y)) = \mathbf{x} + \epsilon \text{sign}(\nabla_{\mathbf{x}} J(\boldsymbol{\theta}, \mathbf{x}, y))$$

$y = "panda"$	$"nematode"$	$"gibbon"$
w/ 57.7%	w/ 8.2%	w/ 99.3 %
confidence	confidence	confidence

(Szegedy *et al.*, Goodfellow *et al.*)

Other Ideas to Reduce Overfitting

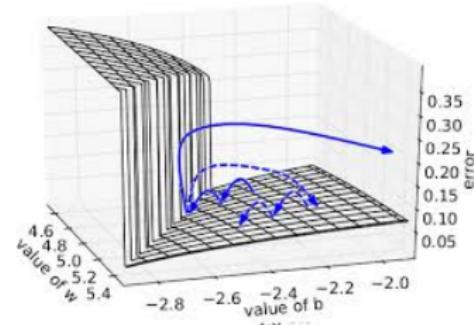
Hard constraints on weights

Gradient Clipping

Inject noise into the system

Enforce sparsity in the neural network

Unsupervised Pre-training



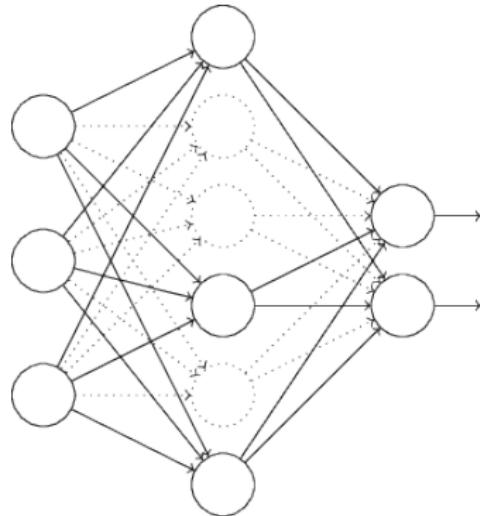
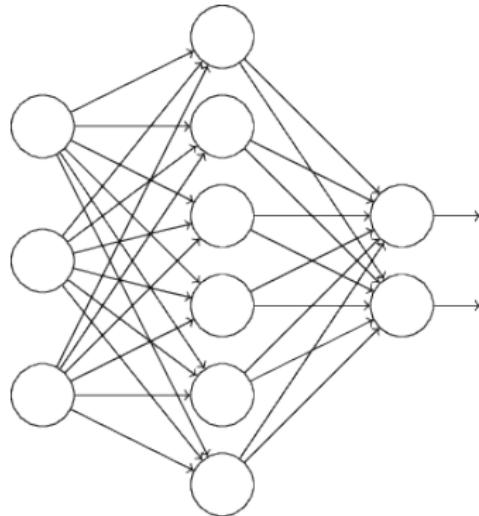
(Bengio *et al.*)

Bagging and Dropout

Bagging (Leo Breiman - 1994)

- ▶ Given dataset $\mathcal{D} = \langle (\mathbf{x}_i, y_i) \rangle_{i=1}^m$, sample $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_k$ of size m from \mathcal{D} with replacement
- ▶ Train classifiers f_1, \dots, f_k on $\mathcal{D}_1, \dots, \mathcal{D}_k$
- ▶ When predicting use majority (or average if using regression)
- ▶ Clearly this approach is not practical for deep networks

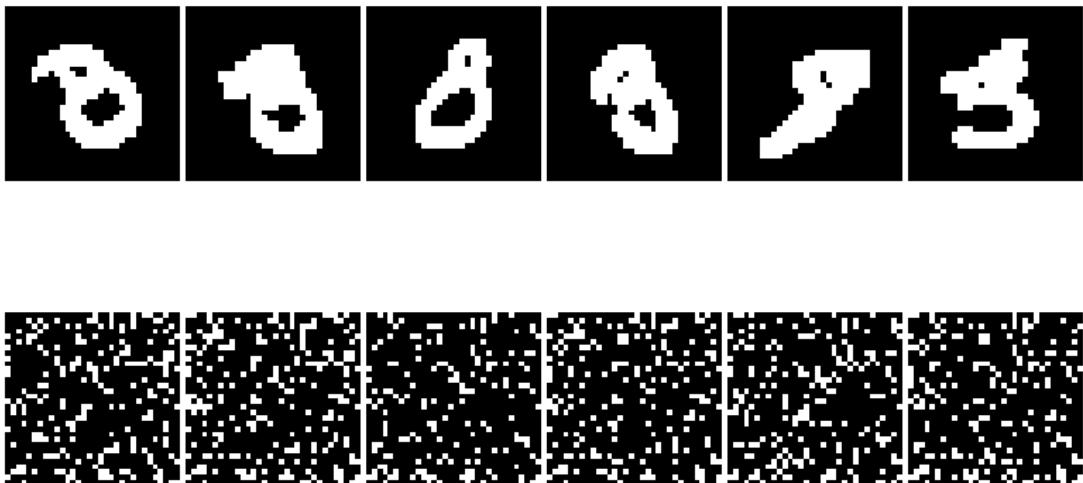
Dropout



- ▶ For input x each hidden unit with probability $1/2$ independently
- ▶ Every input, will have a potentially different mask
- ▶ Potentially exponentially different models, but have “same weights”
- ▶ After training whole network is used by halving all the weights

(Srivastava, Hinton, Krizhevsky, 2014)

Errors Made by MLP for Digit Recognition

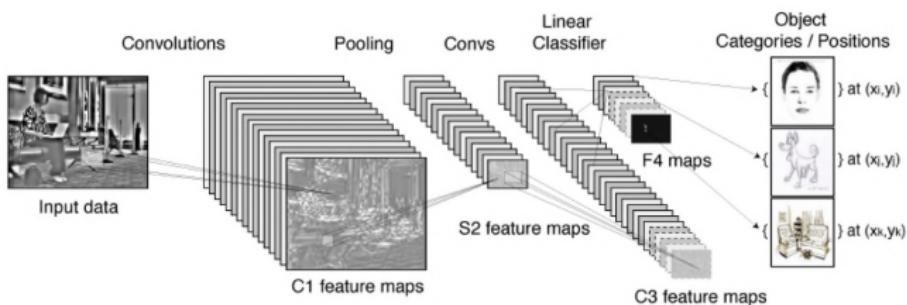


Avoiding Overfitting

- ▶ Use weight sharing a.k.a tied weights in the model
- ▶ Exploit invariances -- translation
- ▶ Exploit locality in images, audio, text, *etc.*
- ▶ Convolutional Neural Networks (convnets)

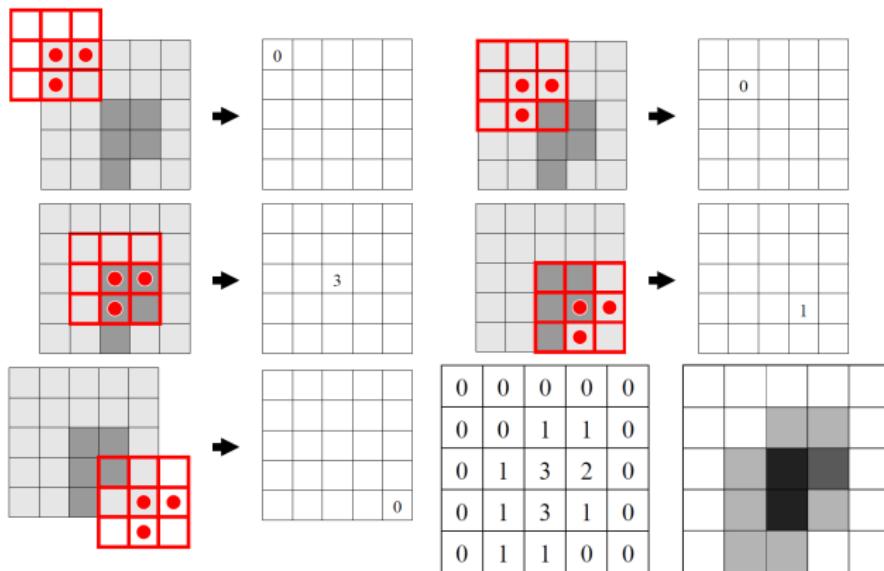
Convolutional Neural Networks (convnets)

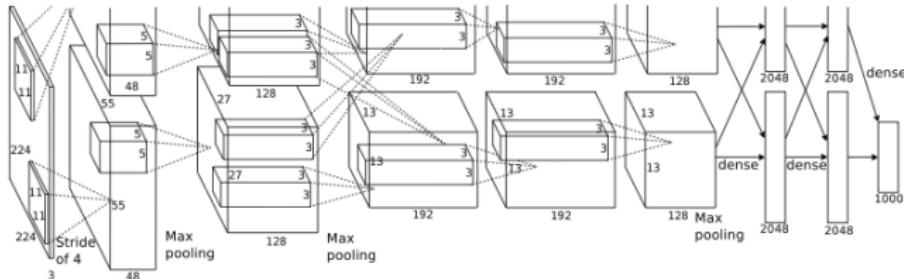
(Fukushima, LeCun, Hinton 1980s)



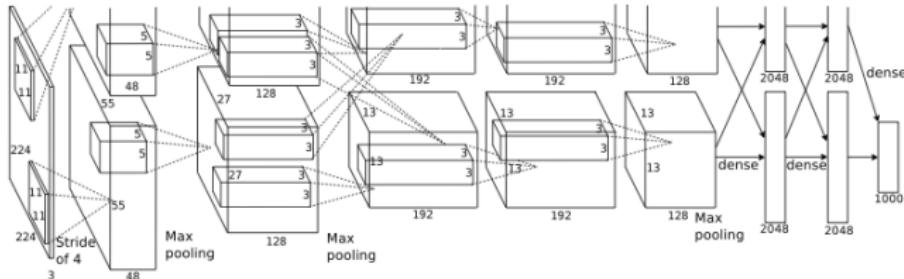
Convolution

Image Convolution

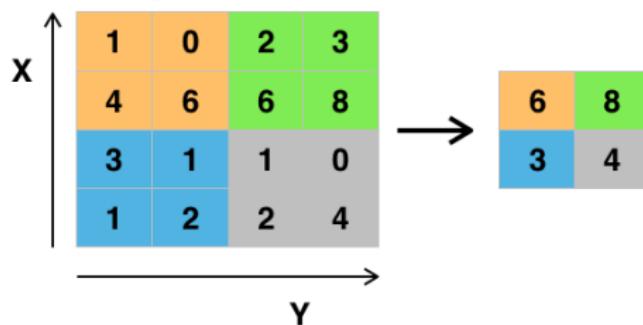




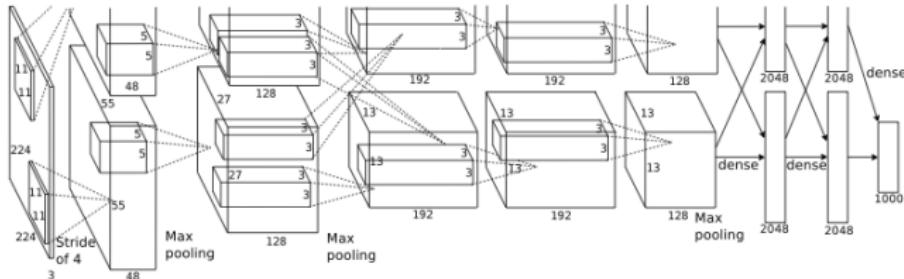
Source: Krizhevsky, Sutskever, Hinton (2012)



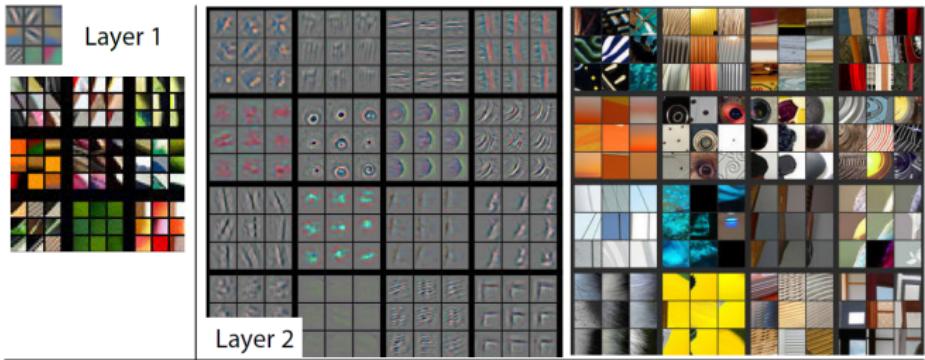
Single depth slice



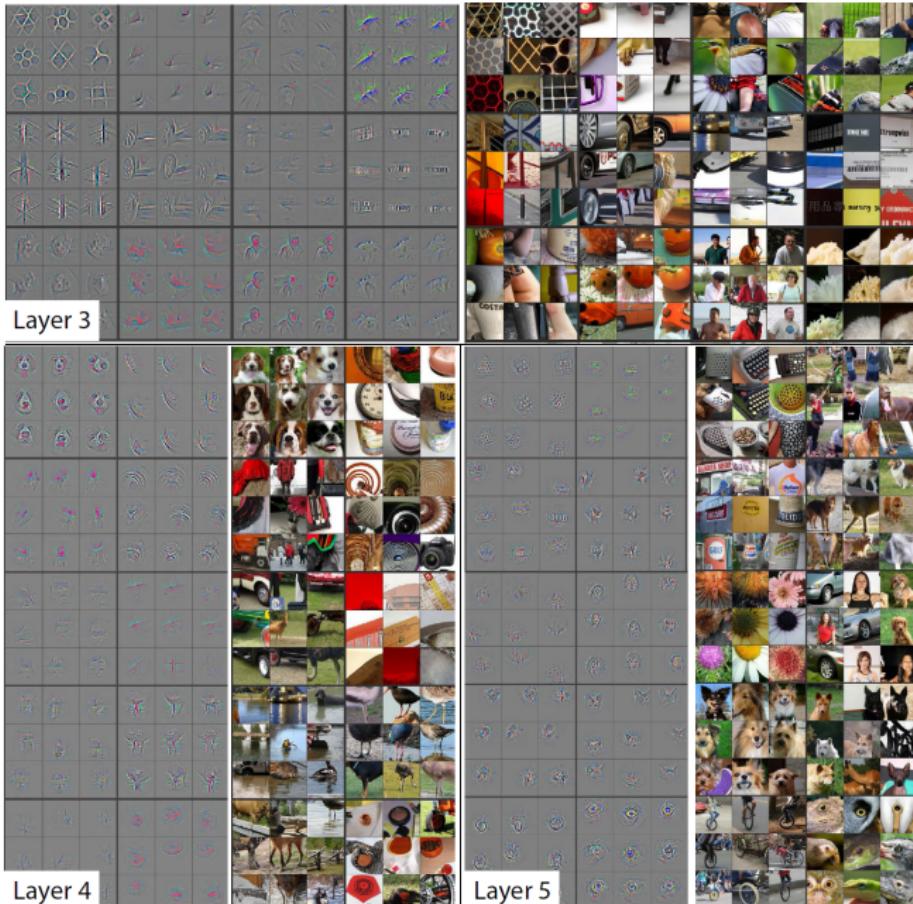
Source: Krizhevsky, Sutskever, Hinton (2012); Wikipedia



Source: Krizhevsky, Sutskever, Hinton (2012)



Source: Zeiler and Fergus (2013)



Source: Zeiler and Fergus (2013)

Convnet in Torch

convnet

```
model = nn.Sequential()
model:add(nn.Reshape(1, 32, 32))
-- layer 1:
model:add(nn.SpatialConvolution(1, 16, 5, 5))
model:add(nn.ReLU())
model:add(nn.SpatialMaxPooling(2, 2, 2, 2))
-- layer 2:
model:add(nn.SpatialConvolution(16, 128, 5, 5))
model:add(nn.ReLU())
model:add(nn.SpatialMaxPooling(2, 2, 2, 2))
-- layer 3:
model:add(nn.Reshape(128*5*5))
model:add(nn.Linear(128*5*5, 200))
model:add(nn.ReLU())
-- output
model:add(nn.Linear(200, 10))
model:add(nn.LogSoftMax())
```

Convolutional Layer

$$z_{i',j',f'}^{l+1} = b_{f'} + \sum_{i=1}^{W_{f'}} \sum_{j=1}^{H_{f'}} \sum_{f=1}^{F_l} a_{i'+i-1,j'+j-1,f}^l w_{i,j,f}^{l+1,f'}$$

$$\frac{\partial z_{i',j',f'}^{l+1}}{\partial w_{i,j,f}^{l+1,f'}} = a_{i'+i-1,j'+j-1,f}^l$$

$$\frac{\partial L}{\partial w_{i,j,f}^{l+1,f'}} = \sum_{i',j'} \delta_{i',j',f'}^{l+1} a_{i'+i-1,j'+j-1,f}^l$$

Convolutional Layer

$$z_{i',j',f'}^{l+1} = b_{f'} + \sum_{i=1}^{W_{f'}} \sum_{j=1}^{H_{f'}} \sum_{f=1}^{F_l} a_{i'+i-1,j'+j-1,f}^l w_{i,j,f}^{l+1,f'}$$

$$\frac{\partial z_{i',j',f'}^{l+1}}{\partial a_{i,j,f}^l} = w_{i-i'+1,j-j'+1,f}^{l+1,f'}$$

$$\frac{\partial L}{\partial a_{i,j,f}^l} = \sum_{i',j',f'} \delta_{i',j',f'}^{l+1} w_{i-i'+1,j-j'+1,f}^{l+1,f'}$$

Max-Pooling Layer

$$b_{i',j'}^{l+1} = \max_{i,j \in \Omega(i',j')} a_{i,j}^l$$

$$\frac{\partial b_{i',j'}^{l+1}}{\partial a_{i,j}^l} = \mathbb{I} \left((i,j) = \operatorname{argmax}_{\tilde{i},\tilde{j} \in \Omega(i',j')} a_{\tilde{i},\tilde{j}}^l \right)$$