

Machine learning - HT 2016

8. Neural Networks

Varun Kanade

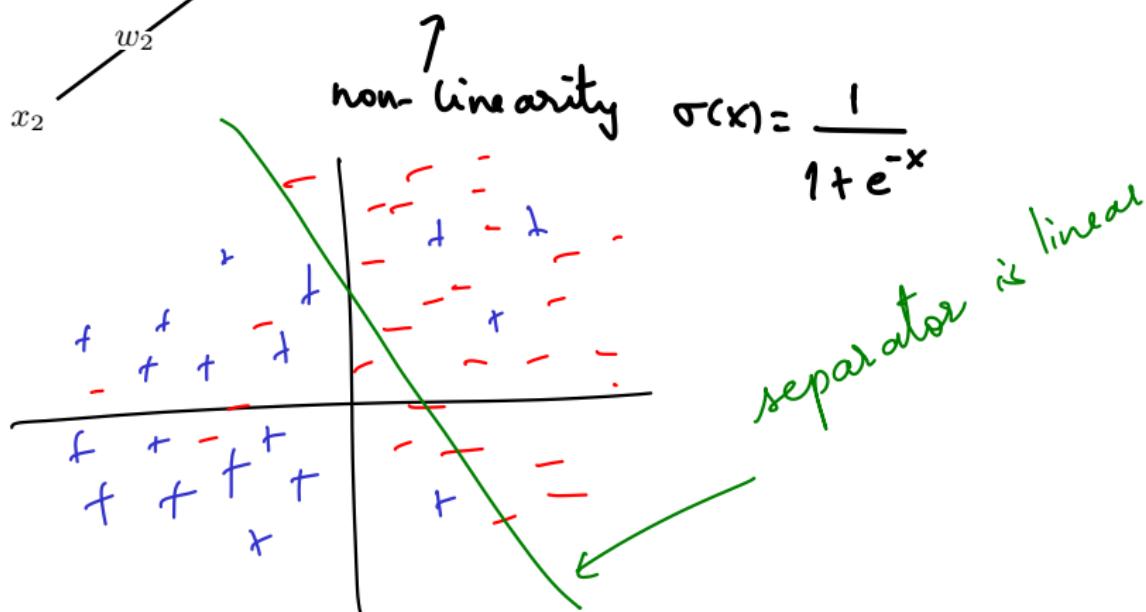
University of Oxford
February 19, 2016

Outline

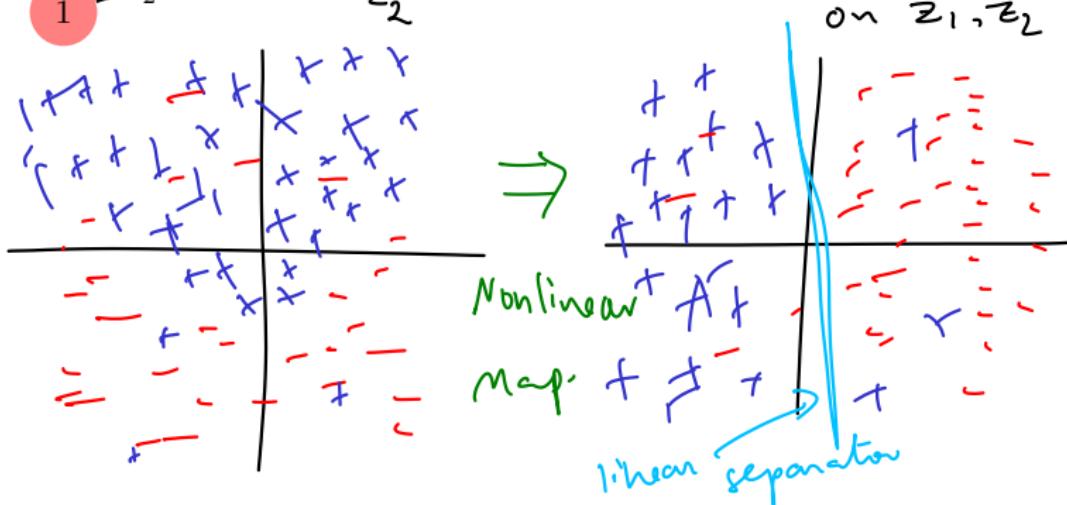
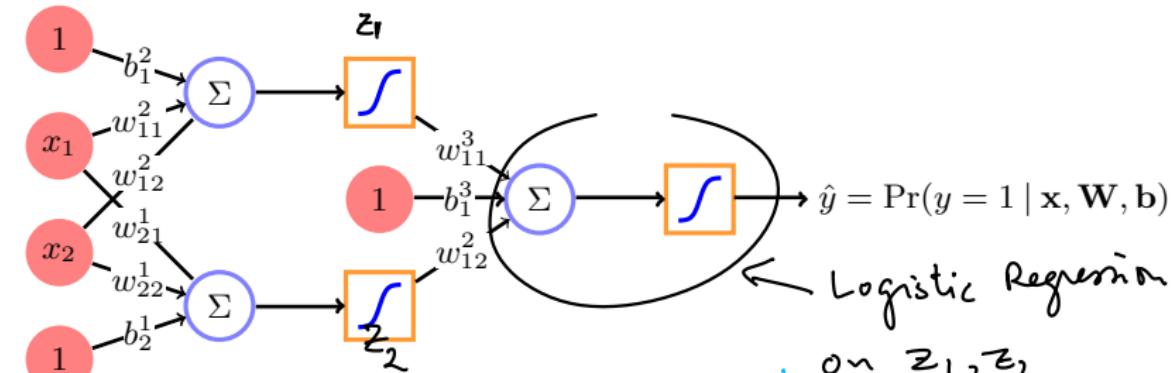
Today, we'll study feedforward neural networks

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

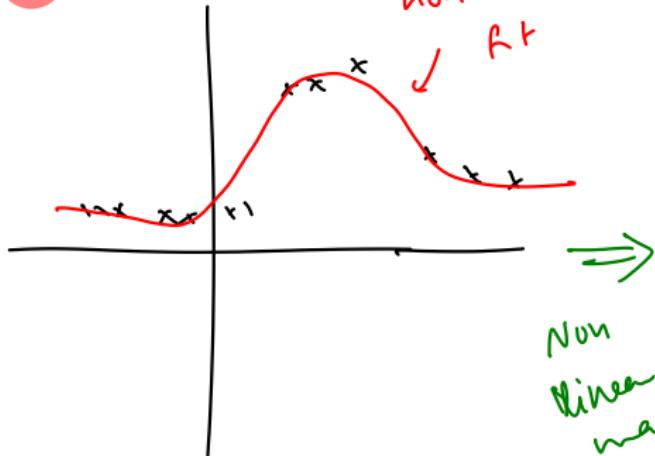
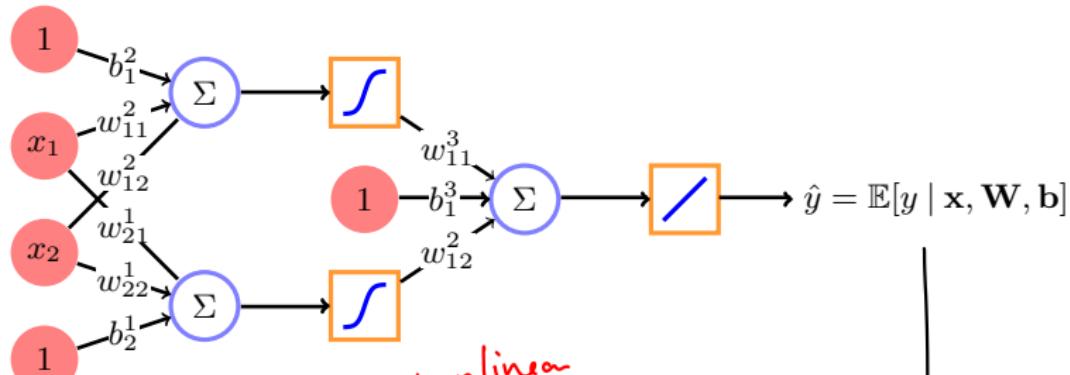
Artificial Neuron : Logistic Regression



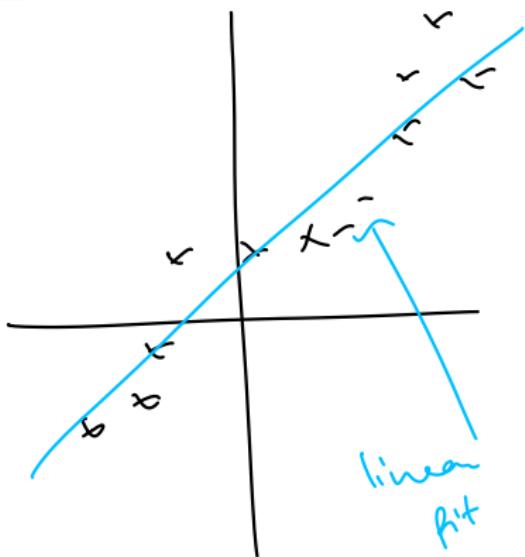
Multilayer Perceptron (MLP) : Classification



Multilayer Perceptron (MLP) : Regression

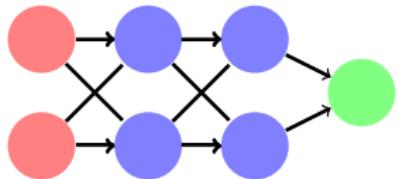


Non
linear
map



linear
fit

A Simple Example



$$z_1^4 = w_{11}^4 \sigma(z_1^3) + w_{12}^4 \sigma(z_2^3) + b_1^4$$

$$z_1^3 = w_{11}^3 \sigma(z_1^2) + w_{12}^3 \sigma(z_2^2) + b_1^3$$

$$z_2^3 = w_{21}^3 \sigma(z_1^2) + w_{22}^3 \sigma(z_2^2) + b_2^3$$

$$z_1^2 = w_{11}^2 x_1 + w_{12}^2 x_2 + b_1^2$$

$$z_2^2 = w_{21}^2 x_1 + w_{22}^2 x_2 + b_2^2$$

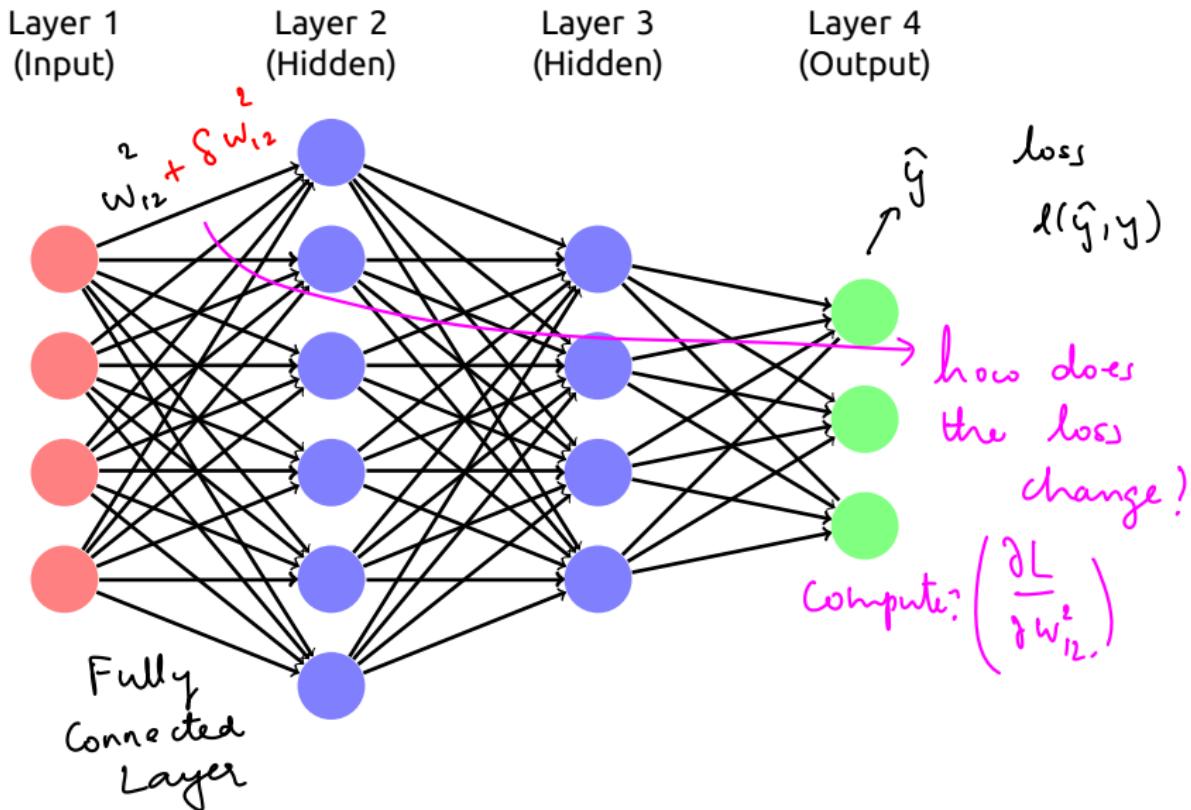
$$\hat{y} = \sigma \left(w_{11}^4 \sigma(w_{11}^3 \sigma(w_{11}^2 x_1 + w_{12}^2 x_2 + b_1^2) + w_{12}^3 \sigma(w_{21}^2 x_1 + w_{22}^2 x_2 + b_2^2) + b_1^3) + w_{12}^4 \sigma(w_{21}^3 \sigma(w_{11}^2 x_1 + w_{12}^2 x_2 + b_1^2) + w_{22}^3 \sigma(w_{21}^2 x_1 + w_{22}^2 x_2 + b_2^2) + b_1^3) + b_1^4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

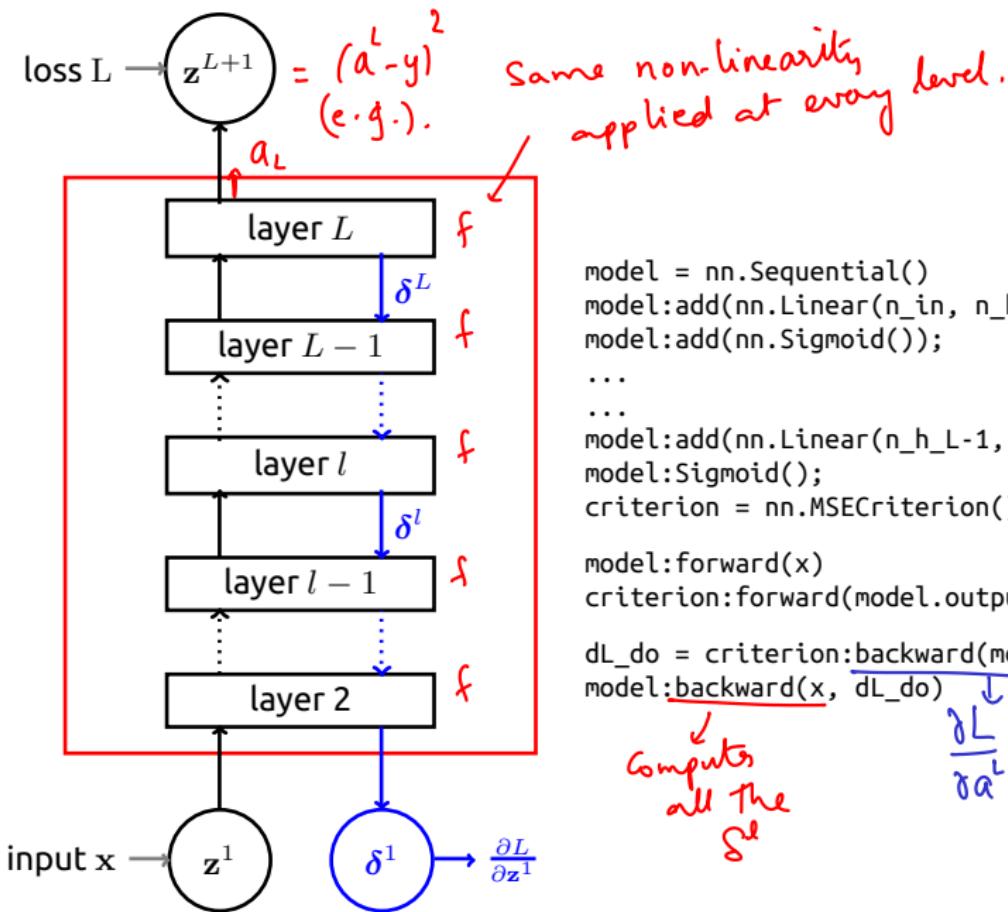
$$L(y, \hat{y}) = \frac{1}{2} (y - \hat{y})^2$$

$$\frac{\partial L}{\partial w_{12}^2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_{12}^2} = (\hat{y} - y) \cdot \left[\sigma'(z_1^4) \cdot \left[w_{11}^4 \cdot \sigma'(z_1^3) \cdot \underbrace{w_{11}^3 \sigma'(z_1^2)}_{+ w_{12}^4 \cdot \sigma'(z_2^3) \cdot \underbrace{w_{21}^3 \sigma'(z_1^2)}_{}} \cdot x_1 \right] \right]$$

Want to exploit common occurrences!

Feedforward Neural Networks





```

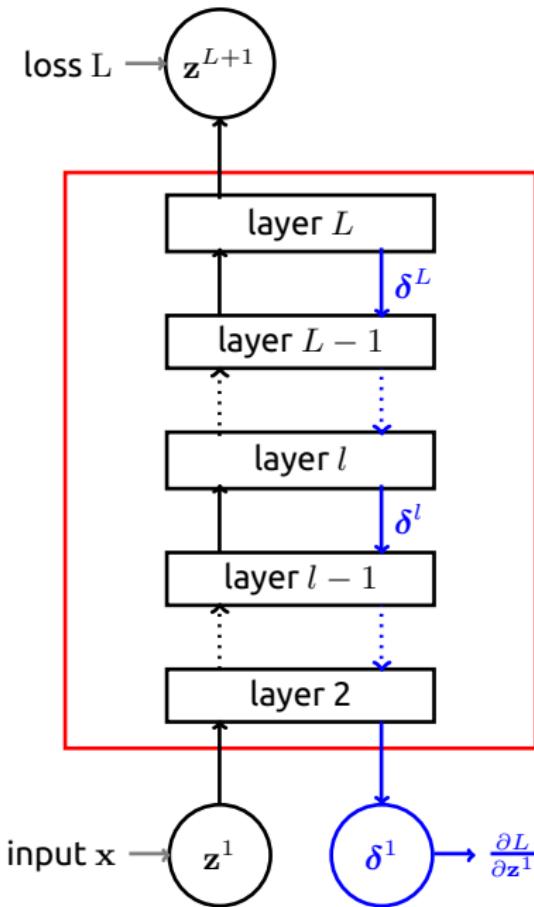
model = nn.Sequential()
model:add(nn.Linear(n_in, n_h1));
model:add(nn.Sigmoid());
...
...
model:add(nn.Linear(n_h_L-1, n_out));
model:Sigmoid();
criterion = nn.MSELoss();

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

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

Computes all the δ^l

$$\frac{\partial L}{\partial a^L}$$



Forward Equations

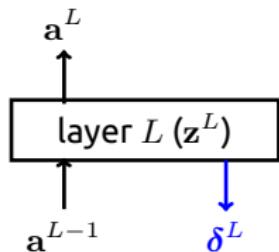
$$(1) \quad \mathbf{z}^1 = \mathbf{x} \text{ (input)}$$

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

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

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

Output Layer



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

$$a^L = f(z^L)$$

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

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

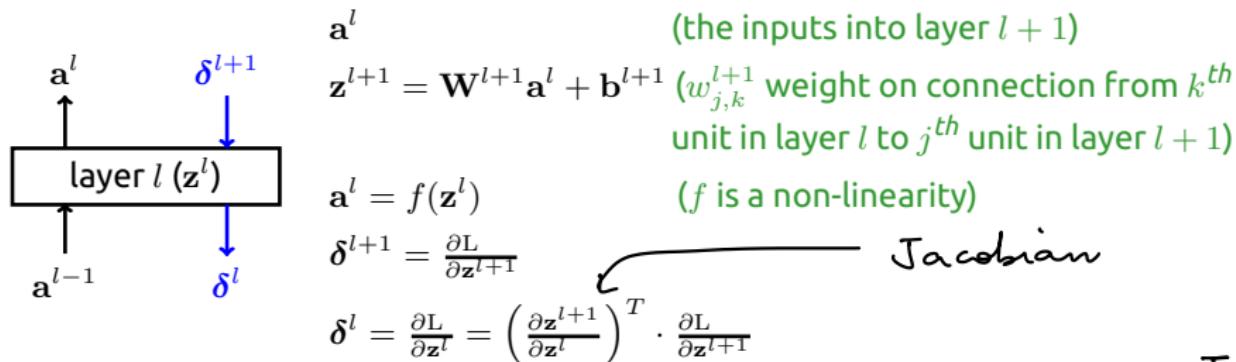
Hadamard
product

(element-wise
vector mult)

$$\frac{\partial L}{\partial z_i^L} = \frac{\partial L}{\partial a_i^L} \cdot f'(z_i^L)$$

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

Back Propagation

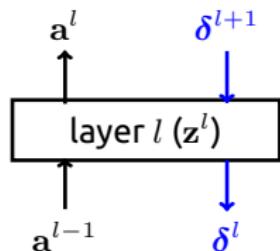


$$\frac{\partial z_i^{l+1}}{\partial z_j^l} = f'(z_j^l) \cdot w_{ij}^{l+1}$$

$$\frac{\partial \mathbf{z}^{l+1}}{\partial \mathbf{z}^l} = \mathbf{W}^{l+1} \cdot \text{Diag}(f'(a_j^l))$$

$$\begin{aligned} \delta^l &= \text{Diag}(f'(a_j^l)) \cdot (\mathbf{W}^{l+1})^T \delta^{l+1} \\ &= f'(a^l) \odot (\mathbf{W}^{l+1})^T \delta^{l+1} \end{aligned}$$

Gradients wrt Parameters



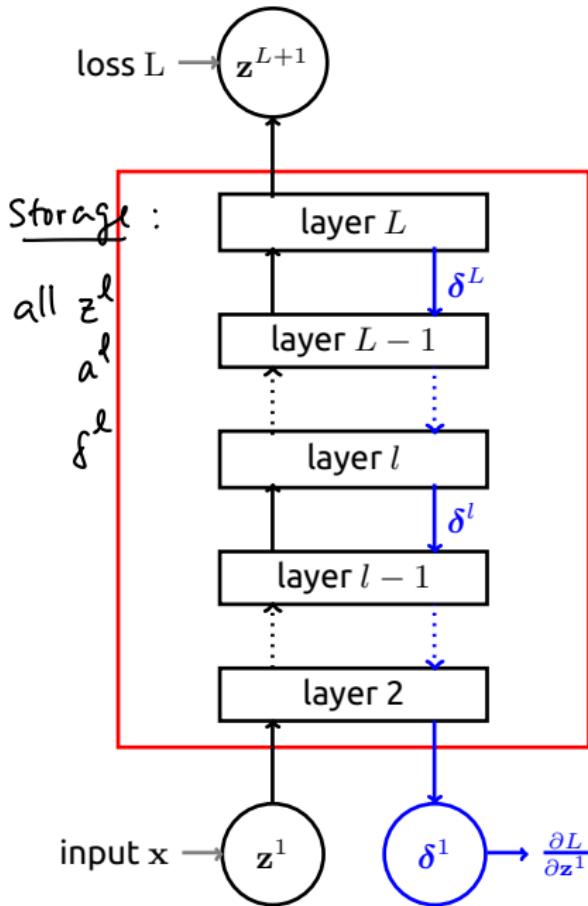
$$z^l = \mathbf{W}^l a^{l-1} + \mathbf{b}^l \quad (\text{$w_{j,k}^l$ weight on connection from k^{th} unit in layer $l-1$ to j^{th} unit in layer l})$$

$$\delta^l = \frac{\partial L}{\partial z^l}$$

$$\begin{aligned} \frac{\partial L}{\partial w_{ij}^l} &= \frac{\partial L}{\partial z_i^l} \cdot a_j^{l-1} \\ &= \delta_i^l \cdot a_j^{l-1} \end{aligned}$$

$$z_i^l = w_{i1}^l \cdot a_1^{l-1} + \dots + w_{in_{l-1}}^l \cdot a_{n_{l-1}}^{l-1} + b_i^l$$

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



Forward Equations

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

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

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

$$(4) \quad L(a^L, y)$$

Save on time, but pay a price in space -

Back-propagation Equations

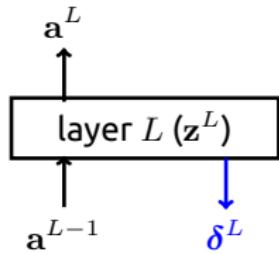
$$(1) \quad \delta^L = \frac{\partial L}{\partial a^L} \odot f'(z^L)$$

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

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

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

Output Layer : LogSoftmax



criterion = nn.CrossEntropyLoss();

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

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

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

$$\frac{\partial \mathcal{L}}{\partial a_y^L} = -1 \quad ; \quad \frac{\partial \mathcal{L}}{\partial a_c^L} = 0 \quad \text{for } c \neq y$$

$$a_i^L = z_i^L - \log(\zeta); \text{ where } \zeta = \sum_j e^{z_j^L}$$

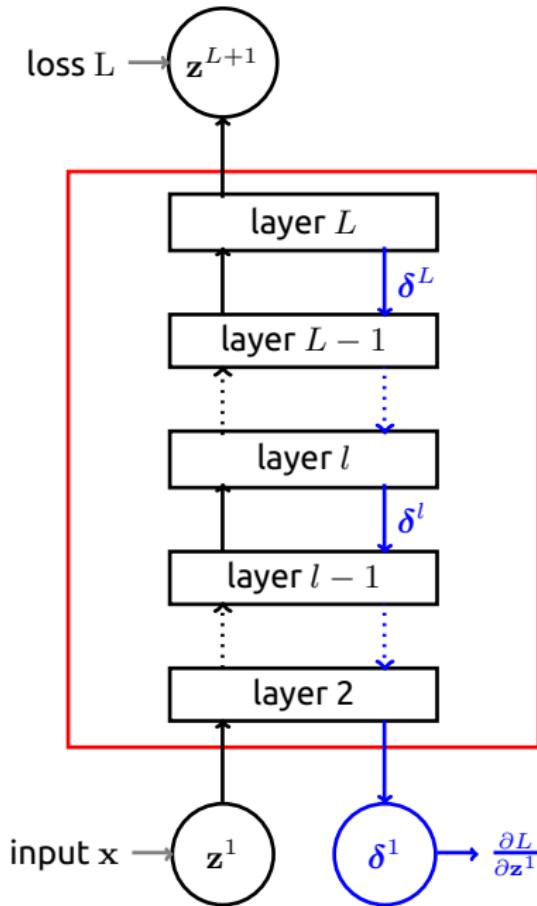
$$\frac{\partial a_i^L}{\partial z_i^L} = 1 - \frac{e^{z_i^L}}{\zeta} = 1 - p_i; \quad p_i = \frac{e^{z_i^L}}{\zeta} \text{ (prob } \hat{y} = i\text{)}$$

$$\frac{\partial a_i^L}{\partial z_j^L} = -\frac{e^{z_j^L}}{\zeta} = -p_j$$

* $\delta_i^L = p_i \quad \forall i \neq y$

$\delta_y^L = p_i - y$

[Remember
multi-class
logistic]



Forward Equations

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

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

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

$$(4) \quad L(a^L, y)$$

Back-propagation Equations

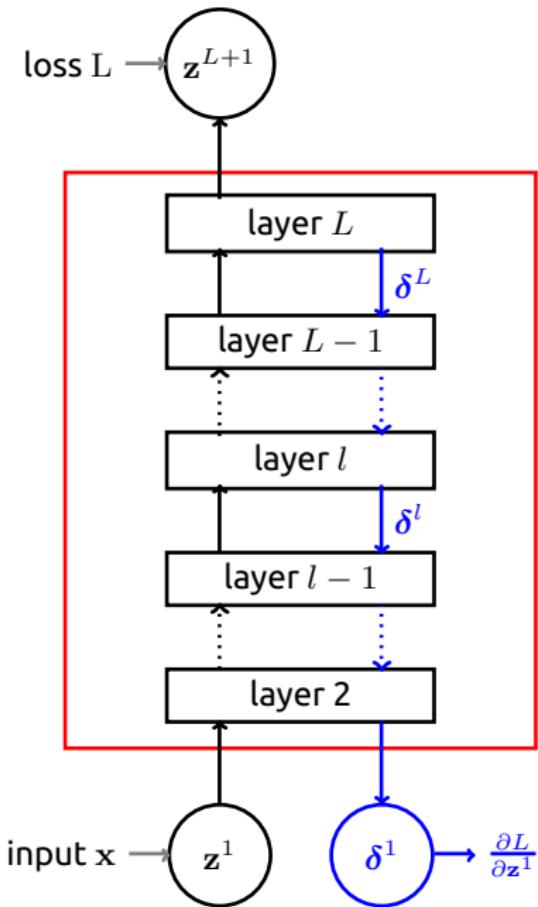
$$(1) \quad \delta^L = \left(\frac{\partial a^L}{\partial z^L} \right)^T \cdot \frac{\partial L}{\partial a^L}$$

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

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

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

If here g^l is some f^l parametrized by W^l, b^l



Forward Equations

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

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

$$(3) \quad \mathbf{a}^l = f(z^l) \quad \leftarrow \text{non-linear}$$

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

pre-activation

activation

Back-propagation Equations

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

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

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

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

Computational Questions

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

As many matrix multiplications as there are linear layers . (Twice—once in the forward & once in the backward pass)
What is the space requirement?

Vectors z^l, a^l, δ^l

Can we process multiple examples together?

Yes minibatching is useful to take advantage of fast tensor operations. But make sure that all parameters can fit in GPU memory.

```
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()
```

Layer 1 Layer 2

```
: 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)
    {
        gradInput : DoubleTensor - size: 4
        weight : DoubleTensor - size: 4x3
        gradWeight : DoubleTensor - size: 4x3
        gradBias : DoubleTensor - size: 4
        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)
    {
        gradInput : DoubleTensor - size: 2
        weight : DoubleTensor - size: 2x4
        gradWeight : DoubleTensor - size: 2x4
        gradBias : DoubleTensor - size: 2
        addBuffer : DoubleTensor - size: 500
        bias : DoubleTensor - size: 2
        output : DoubleTensor - size: 500x2
    }
    4:
        nn.Sigmoid
    {
        gradInput : DoubleTensor - size: 500x2
        output : DoubleTensor - size: 500x2
    }
}

```

δ^1
 $\frac{\partial L}{\partial a^1}$
 δ^2

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?

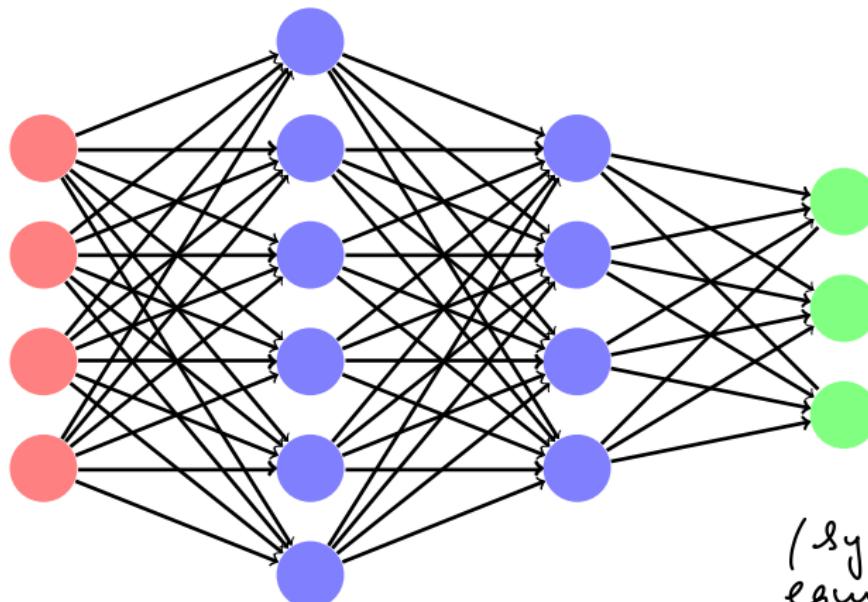
Training Feedforward Deep Networks

Layer 1
(Input)

Layer 2
(Hidden)

Layer 3
(Hidden)

Layer 4
(Output)

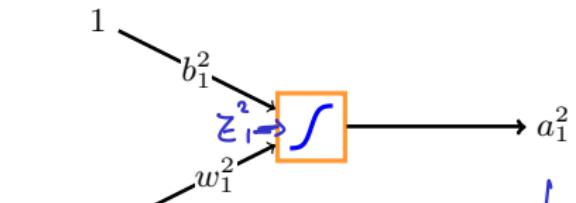


(Symmetries lead to equivalent weights (eg any permutation of hidden units.) We are more

worried about BAD local minima.

Why do we get non-convex optimisation problem?

A toy example



Target is $y = \frac{1-x}{2}$

$$x = +1 \\ n = -1$$

$$y = 0 \\ y = 1$$

$$z_1^2 = w_1^2 x + b_1^2$$

$$L = \frac{1}{2} (y - \sigma(z_1^2))^2$$

$$\frac{\partial L}{\partial z_1^2} = (y - \sigma(z_1^2)) \sigma'(z_1^2)$$

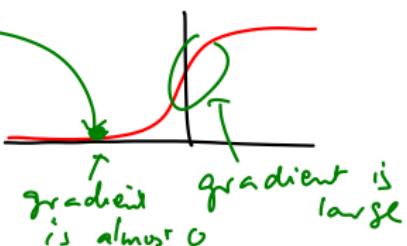
CROSS ENTROPY

$$L = -y \log(\sigma(z_1^2)) + (1-y) \log(1-\sigma(z_1^2))$$

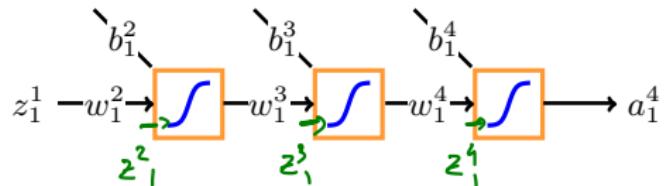
Say $x = -1, y = 1,$

then $\frac{\partial L}{\partial z_1^2} = 1 - \sigma(z_1^2)$

$$\left| \begin{array}{l} x = -1 \\ w_1^2 = 2, b_1^2 = 0.5 \\ z_1^2 = -2.5 \end{array} \right.$$



Propagating Gradients Backwards



(SATURATION)

$$z^4_1 = \omega_1^4 \sigma(z^3_1) + b^4_1 ; L : \text{CROSS ENTROPY.}$$

$$\frac{\partial L}{\partial z^4_1} = y - \sigma(z^4_1)$$

$$\sigma' \in (0, 1/4)$$

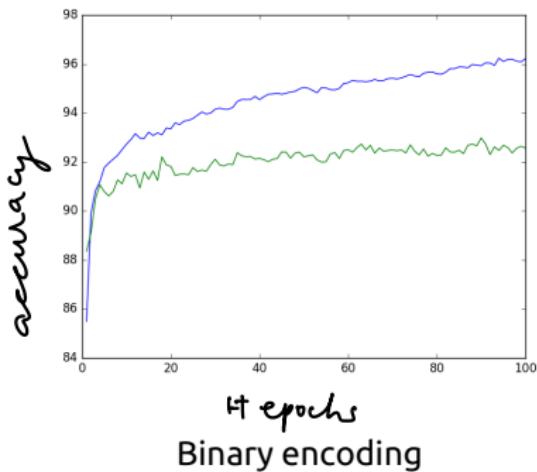
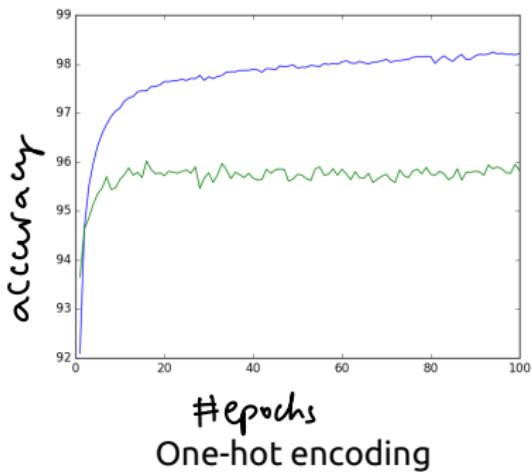
$$\frac{\partial L}{\partial z^3_1} = \frac{\partial L}{\partial z^4_1} \cdot (\omega_1^4 \cdot \sigma'(z^3_1))$$

$$\frac{\partial L}{\partial z^2_1} = \frac{\partial L}{\partial z^4_1} \cdot (\omega_1^4 \cdot \sigma'(z^3_1))(\omega_1^3 \cdot \sigma'(z^2_1))$$

(VANISHING GRADIENT)

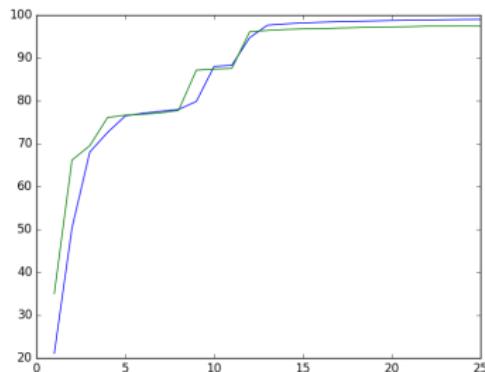
← gradient
across many layers can
go to 0

Digit Classification Problem on Sheet 4

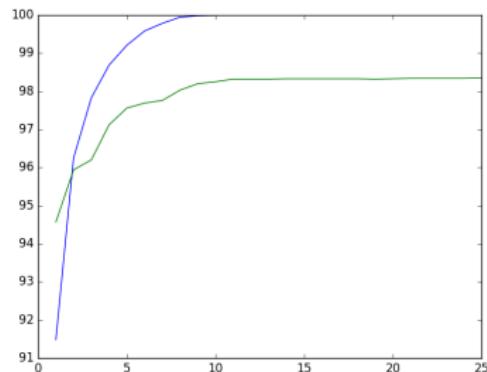


Digit Classification: Squared Loss vs Cross Entropy

uses 1024 hidden units



Squared Error



Cross Entropy

(See paper by Glorot and Bengio (2010))

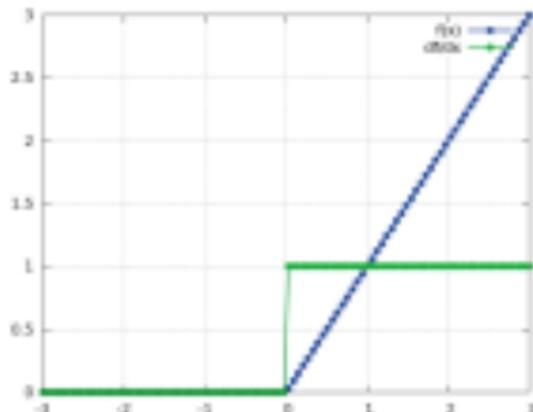
Avoiding Saturation

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

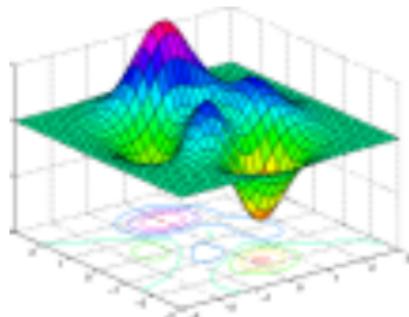
Why is initialising important?

Non-convexity. Initialization can affect which local minima is reached.

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

small values. If $z = \sum_{i=1}^n w_i a_i$

What if it were a ReLU unit? e.g., choose $w_i \in [-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}]$ at random.



How about the biases?

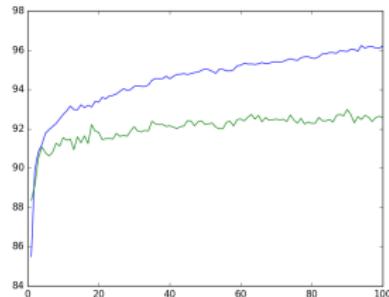
Avoiding Overfitting

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?

- (i) maintain and compute validation error
- (ii) store model when necessary

What are the advantages?

- (i) efficient, compared to other methods

Early stopping leads to better generalisation

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

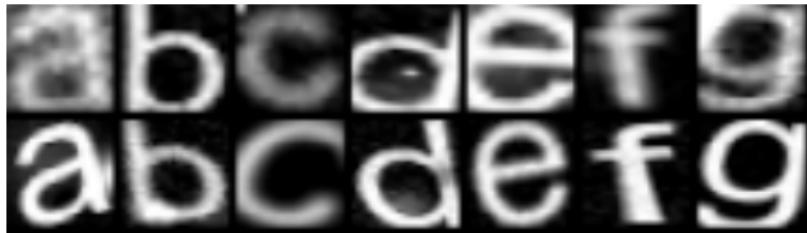
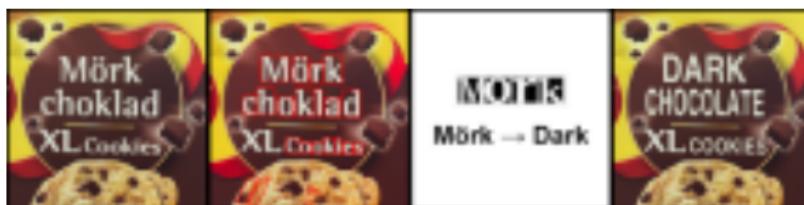
Add Data

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

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

x $\text{sign}(\nabla_x J(\theta, x, y))$ $x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$
 $y = \text{"panda"}$ "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

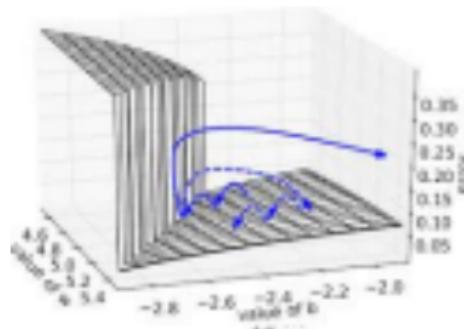
Hard constraints on weights

Gradient Clipping

Inject noise into the system

Enforce sparsity in the neural network

Unsupervised Pre-training



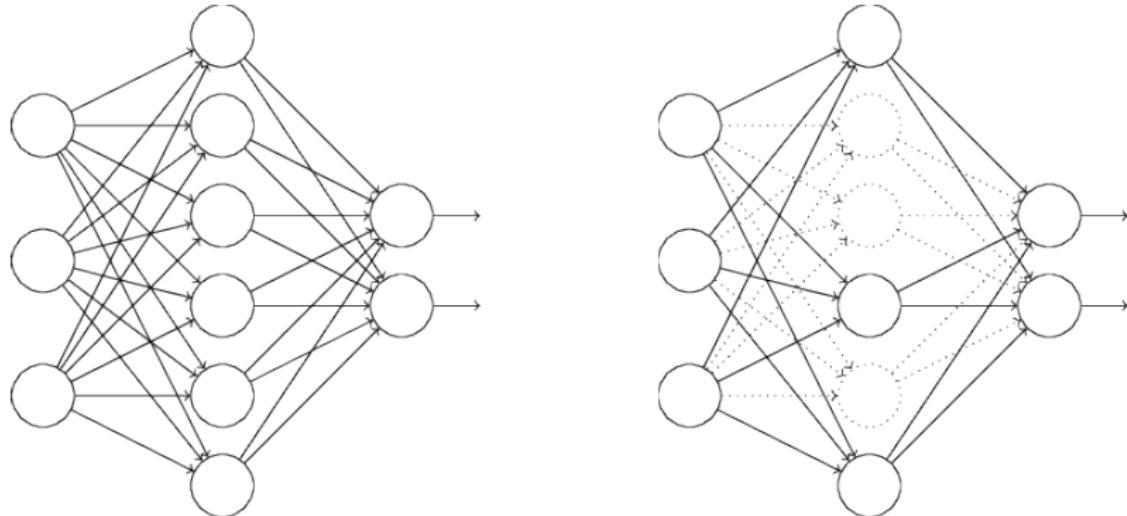
(Bengio *et al.*)

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

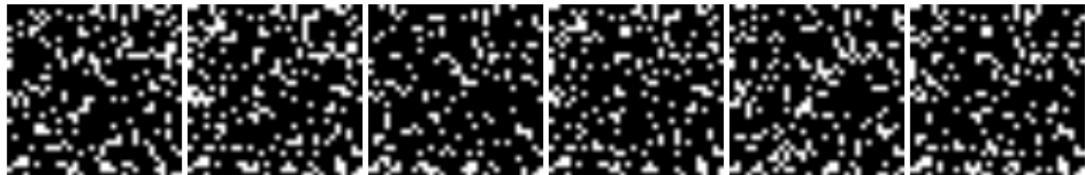
Dropout



- ▶ 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

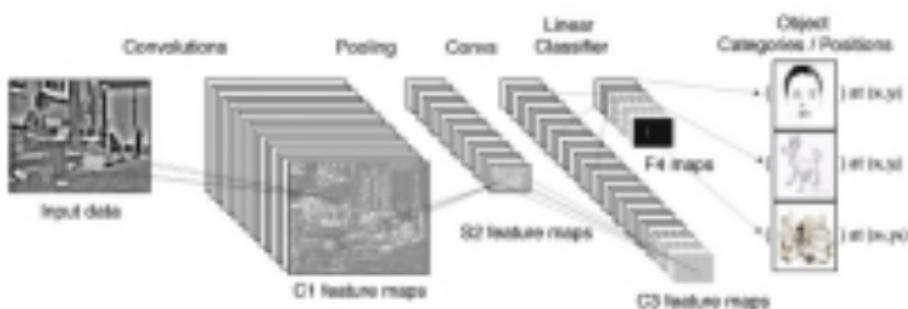


Avoiding Overfitting

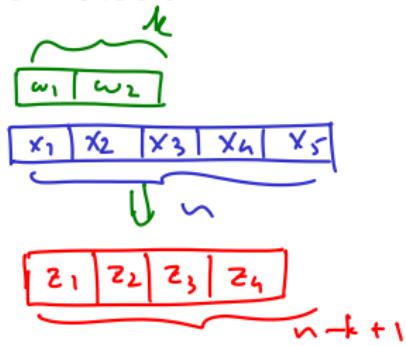
- ▶ Use weight sharing a.k.a tied weights in the model
- ▶ Exploit invariances -- translation
- ▶ Exploit locality in images, audio, text, *etc.*
- ▶ Convolutional Neural Networks (convnets)

Convolutional Neural Networks (convnets)

(Fukushima, LeCun, Hinton 1980s)



Convolution



$$z_1 = x_1 w_1 + x_2 w_2$$

$$z_2 = x_2 w_1 + x_3 w_2$$

$$z_3 = x_3 w_1 + x_4 w_2$$

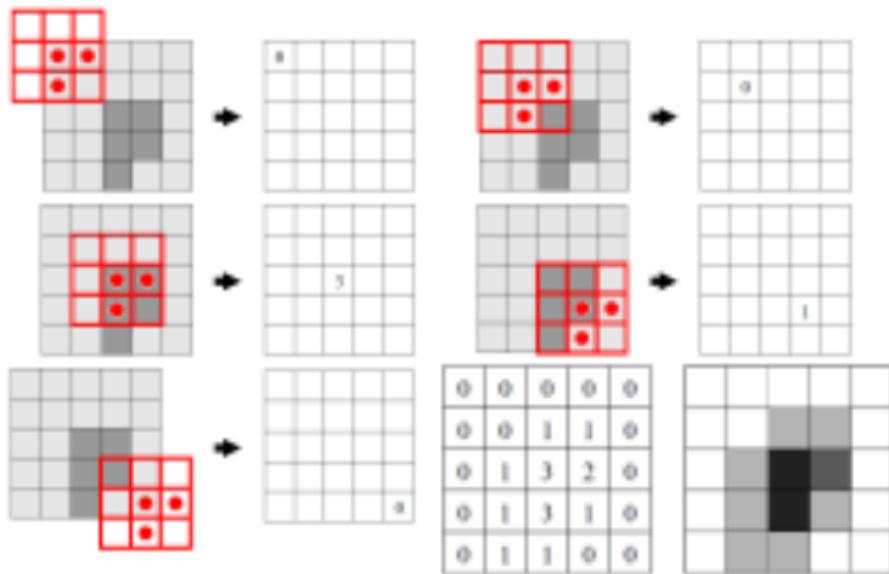
$$z_4 = x_4 w_1 + x_5 w_2$$

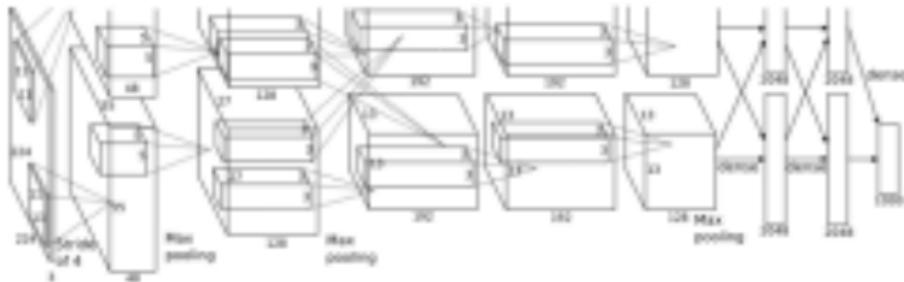
zero padding to get suitable length.

stride: don't take every consecutive position.

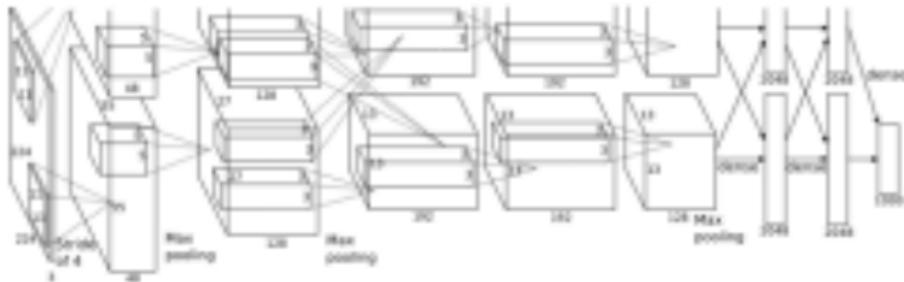
[Mathematically: convolution is defined somewhat differently.
This is actually correlation.]
[See: kernel flipping]

Image Convolution

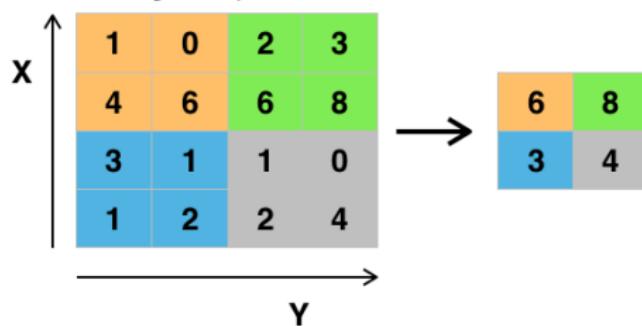




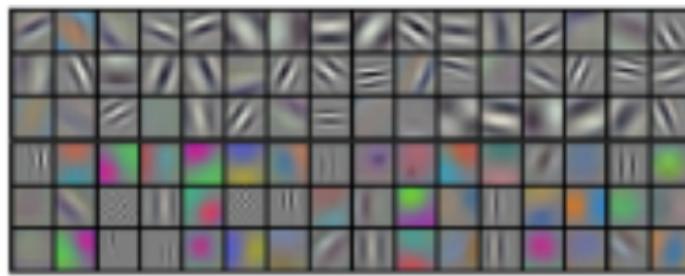
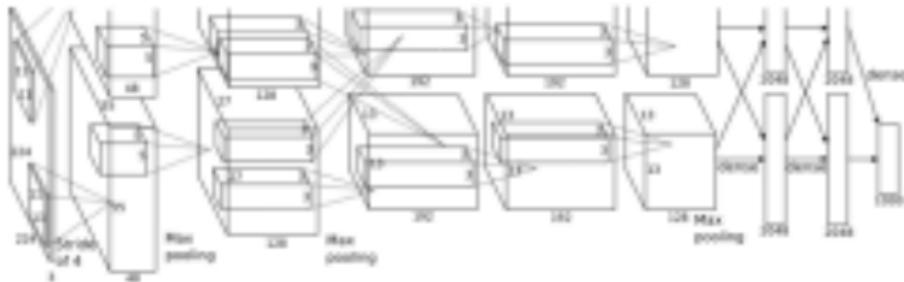
Source: Krizhevsky, Sutskever, Hinton (2012)



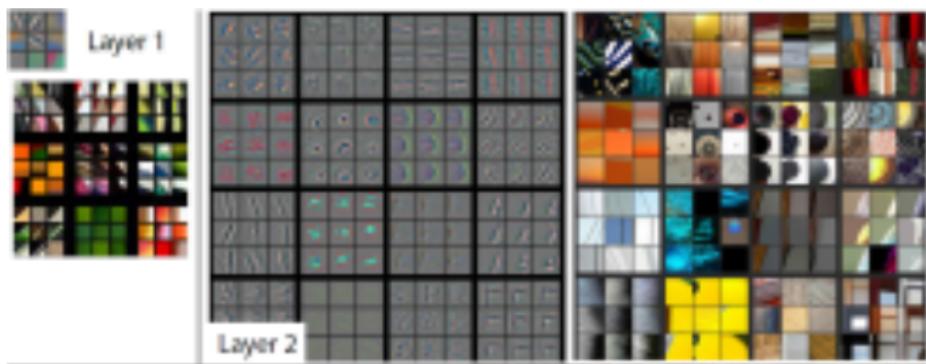
Single depth slice



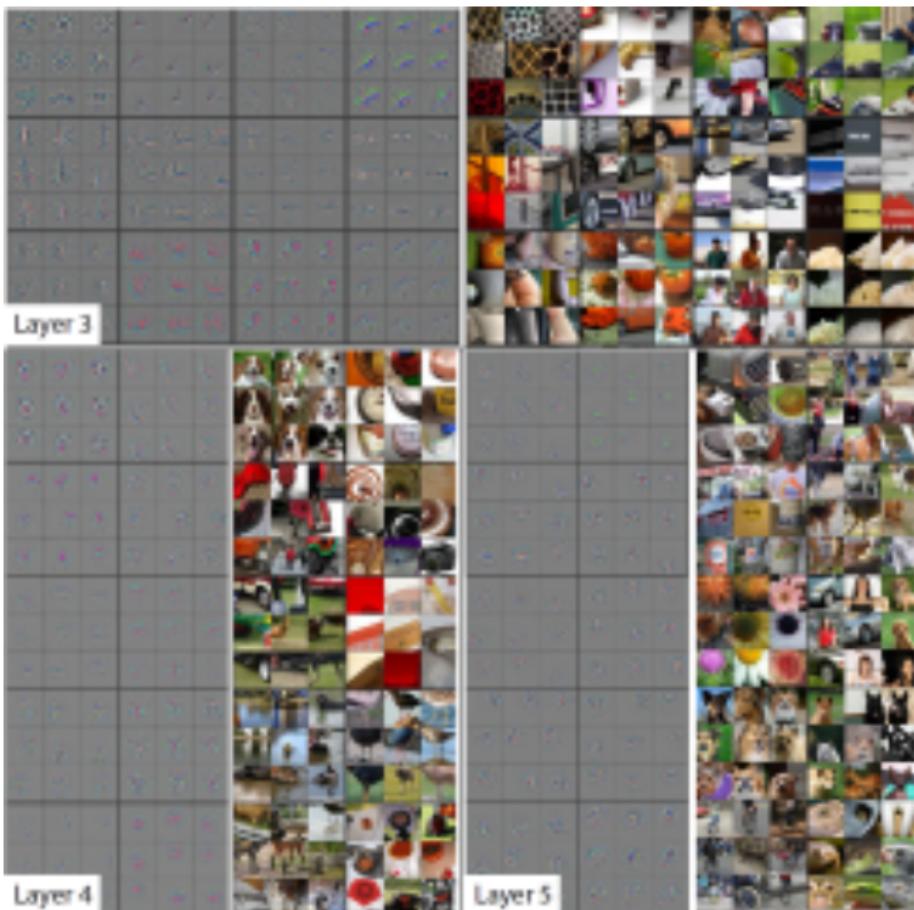
Source: Krizhevsky, Sutskever, Hinton (2012); Wikipedia



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

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

← Domain over which pooling is performed.

{The derivative can be undefined when
argmax is actually a set, and not a single
point}