

# 4F10: Deep Learning

Mark Gales

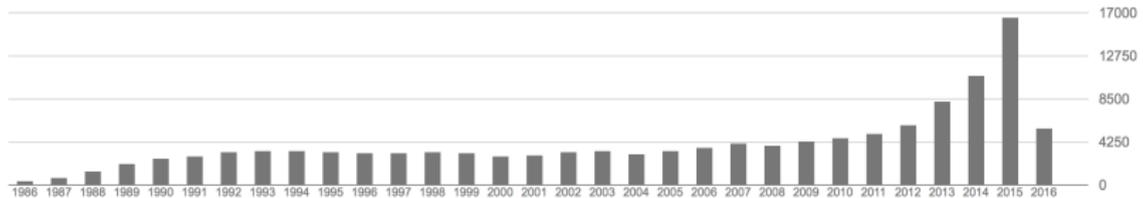
Michaelmas 2016

# What is Deep Learning?

From Wikipedia:

*Deep learning is a branch of machine learning based on a set of algorithms that attempt to model high-level abstractions in data by using multiple processing layers, with complex structures or otherwise, composed of multiple non-linear transformations.*

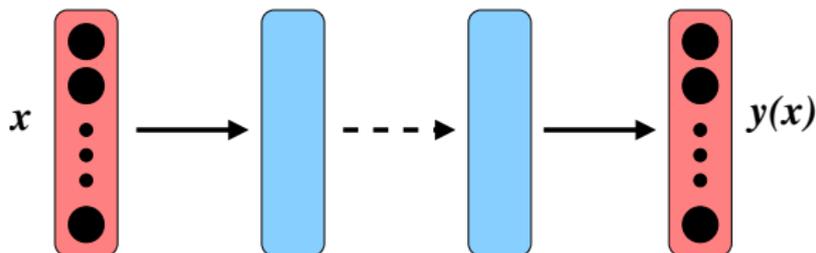
# The Rise of Deep Learning (June 2016)



- Plot shows citations to [Geoff Hinton](#) papers
  - highly influential researcher in deep learning

- Basic Building Blocks
  - neural network architectures
  - activation functions
- Error Back Propagation
  - single-layer perceptron (motivation)
  - multiple-layer perceptron
- Optimisation
  - gradient descent refinement
  - second-order approaches, and use of curvature
  - initialisation
- Example
  - encoder-decoder models for sequence-to-sequence models
  - and attention for sequence-to-sequence modelling

# Basic Building Blocks

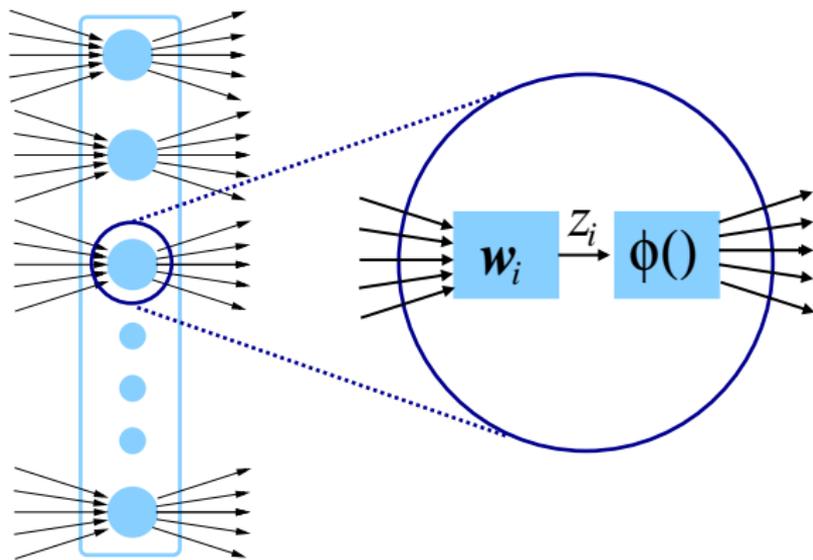


- General mapping process from input  $x$  to output  $y(x)$

$$y(x) = \mathcal{F}(x)$$

- deep refers to number of **hidden** layers
- Output from the previous layer connected to following layer:
  - $x^{(k)}$  is the input to layer  $k$
  - $x^{(k+1)} = y^{(k)}$  the output from layer  $k$

# Neural Network Layer/Node



- General form for layer  $k$ :

$$y_i^{(k)} = \phi(\mathbf{w}'_i \mathbf{x}^{(k)} + b_i) = \phi(z_i^{(k)})$$

# Initial Neural Network Design Options

- The input and outputs to the network are defined
  - able to select number of **hidden** layers
  - able to select number of **nodes** per hidden layer
- Increasing layers/nodes increases model parameters
  - need to consider how well the network **generalises**
- For **fully connected** networks, number of parameters ( $N$ ) is

$$N = d \times N^{(1)} + K \times N^{(L)} + \sum_{k=1}^{L-1} N^{(k)} \times N^{(k+1)}$$

- $L$  is the number of **hidden layers**
  - $N^{(k)}$  is the number of nodes for layer  $k$
  - $d$  is the input vector size,  $K$  is the output size
- Designing “good” networks is complicated ...

- Heaviside (or step/threshold) function: output binary

$$\phi(z_i) = \begin{cases} 0, & z_i < 0 \\ 1, & z_i \geq 0 \end{cases}$$

- Sigmoid function: output continuous,  $0 \leq y_i(\mathbf{x}) \leq 1$ .

$$\phi(z_i) = \frac{1}{1 + \exp(-z_i)}$$

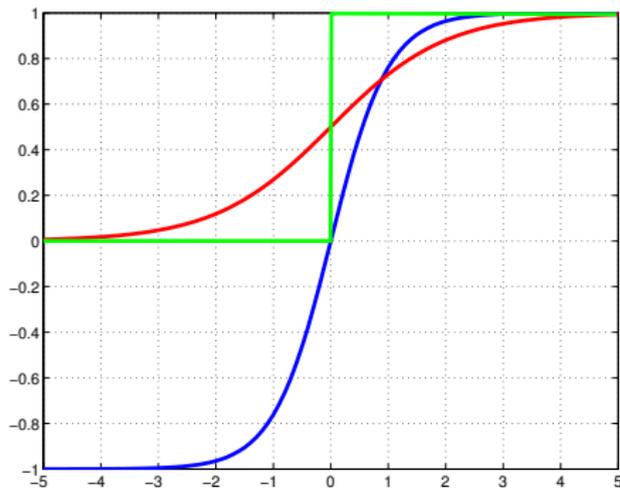
- Softmax function: output  $0 \leq y_i(\mathbf{x}) \leq 1$ ,  $\sum_{i=1}^n y_i(\mathbf{x}) = 1$ .

$$\phi(z_i) = \frac{\exp(z_i)}{\sum_{j=1}^n \exp(z_j)}$$

- Hyperbolic tan function: output continuous,  $-1 \leq y_i(\mathbf{x}) \leq 1$ .

$$\phi(z_i) = \frac{\exp(z_i) - \exp(-z_i)}{\exp(z_i) + \exp(-z_i)}$$

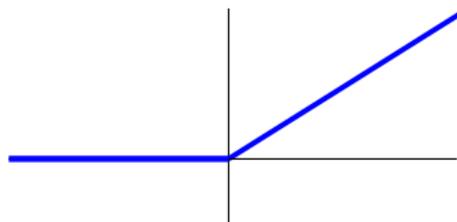
# Activation Functions



- Activation functions:
  - step function (green)
  - sigmoid function (red)
  - tanh function (blue)

- softmax, usual output layer for classification tasks
- sigmoid/tanh, often used for hidden layers

- Alternative activation function: Rectified Linear Units



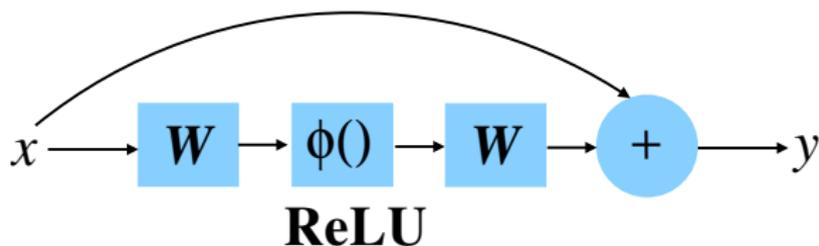
$$\phi(z_i) = \max(0, z_i)$$

- Related activation function **noisy ReLU**:

$$\phi(z_i) = \max(0, z_i + \epsilon); \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

- efficient, no exponential/division, rapid convergence in training
- Leaky ReLU** also possible

$$\phi(z_i) = \begin{cases} z_i; & z_i \geq 0; \\ \alpha z_i & z_i < 0 \end{cases}$$

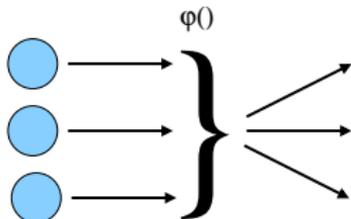


- Modify layer to model the **residual**

$$\mathbf{y}(\mathbf{x}) = \mathcal{F}(\mathbf{x}) + \mathbf{x}$$

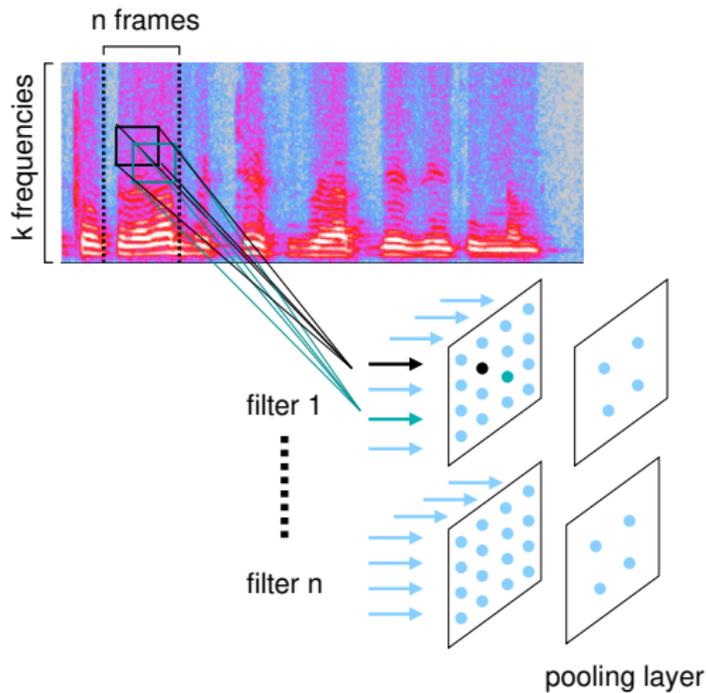
- allows deeper networks to be built
  - **deep residual learning**
- Links to **highway connections**

- Possible to **pool** the output of a set of node
  - reduces the number of weights to connect layers together



- A range of functions have been examined
  - **maxout**  $\phi(y_1, y_2, y_3) = \max(y_1, y_2, y_3)$
  - **soft-maxout**  $\phi(y_1, y_2, y_3) = \log(\sum_{i=1}^3 \exp(y_i))$
  - **p-norm**  $\phi(y_1, y_2, y_3) = (\sum_{i=1}^3 |y_i|)^{1/p}$
- Has also been applied for unsupervised adaptation

# Convolutional Neural Networks [14, 1]



- Various parameters control the form of the CNN:
  - **number** (depth): how many filters to use
  - **receptive field** (filter size): height/width/depth ( $h \times w \times d$ )
  - **stride**: how far filter moves in the convolution
  - **dilation**: “gaps” between filter elements
  - **zero-padding**: do you pad the edges of the “image” with zeroes
- Filter output can be stacked to yield depth for next layer
- To illustrate the impact consider 1-dimensional case - default:
  - **zero-padding**, **stride=1**, **dilation=0**

sequence: 1, 2, 3, 4, 5, 6      filter: 1, 1, 1

default: 3, 6, 9, 12, 15, 11      no padding: 6, 9, 12, 15

dilation=1: 4, 6, 9, 12, 8, 10      stride=2: 3, 9, 15

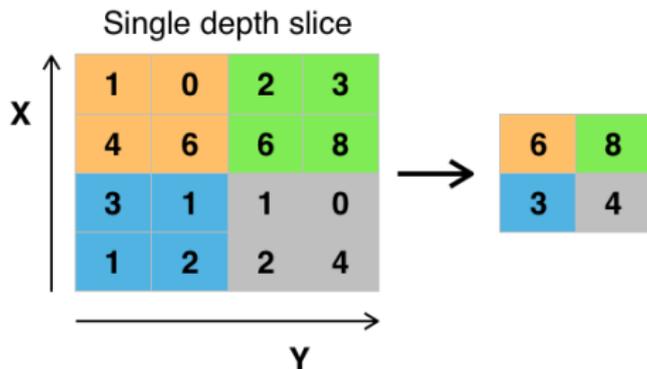
# Convolutional Neural Network (cont)

- For a 2-dimensional image the convolution can be written as

$$\phi(z_{ij}) = \phi \left( \sum_{kl} w_{kl} x_{(i-k)(j-l)} \right)$$

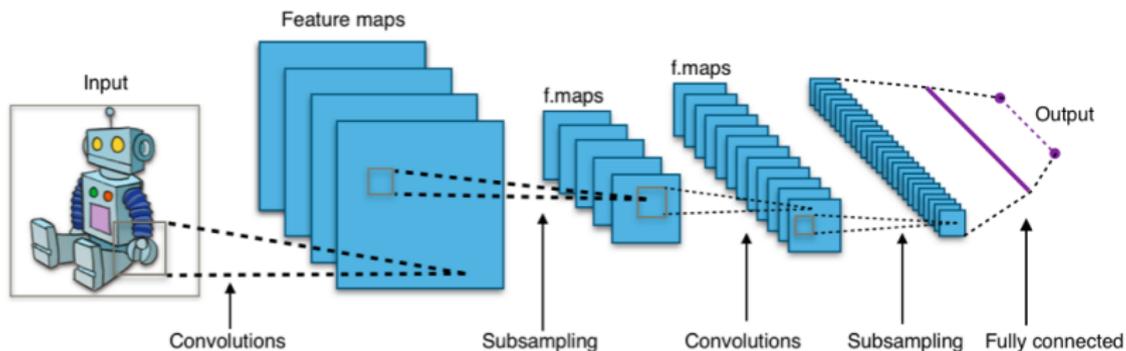
- $x_{ij}$  is the image value at point  $i, j$
- $w_{kl}$  is the weight at point  $k, l$
- $\phi$  is the non-linear activation function
- For a  $5 \times 5$  receptive field (no dilation)
  - $k \in \{-2, -1, 0, 1, 2\}$
  - $l \in \{-2, -1, 0, 1, 2\}$
- Stride determines how  $i$  and  $j$  vary as we move over the image

# CNN Max-Pooling (Subsampling)



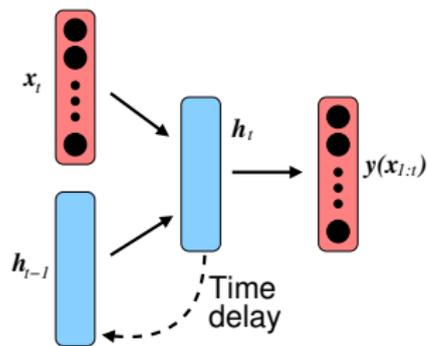
- **Max-Pooling** over a  $2 \times 2$  filter on  $4 \times 4$  “image”
  - **stride** of 2 - yields output of  $2 \times 2$
- Possible to also operate with a stride of 1 **overlapping pooling**

# Simple CNN Example



- Simple five layer classifier for images
  - two convolutional layers each followed by
  - pooling layers (two)
  - with a fully connected network and softmax activation function

- Consider a **causal** sequence of observations  $\mathbf{x}_{1:t} = \{\mathbf{x}_1, \dots, \mathbf{x}_t\}$



- Introduce recurrent units

$$\mathbf{h}_t = \mathbf{f}^h(\mathbf{W}_h^f \mathbf{x}_t + \mathbf{W}_h^r \mathbf{h}_{t-1} + \mathbf{b}_h)$$

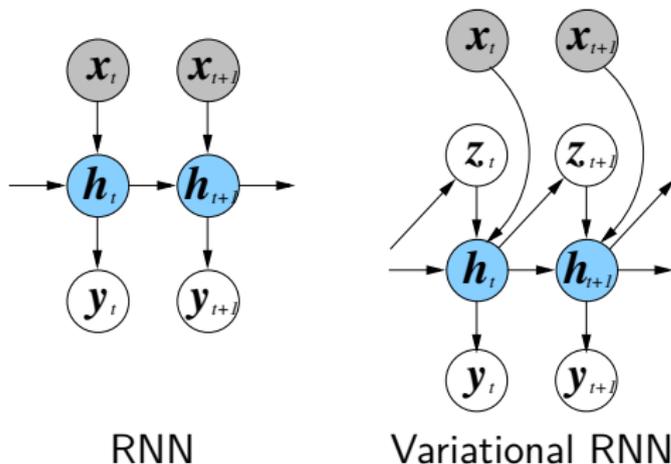
$$\mathbf{y}(\mathbf{x}_{1:t}) = \mathbf{f}^f(\mathbf{W}_y \mathbf{h}_t + \mathbf{b}_y)$$

- $\mathbf{h}_t$  **history vector** at time  $t$
- Two history weight matrices
  - $\mathbf{W}_h^f$  forward,  $\mathbf{W}_h^r$  recursion
- Uses approximation to model **history of observations**

$$\mathcal{F}(\mathbf{x}_{1:t}) = \mathcal{F}(\mathbf{x}_t, \mathbf{x}_{1:t-1}) \approx \mathcal{F}(\mathbf{x}_t, \mathbf{h}_{t-1}) \approx \mathcal{F}(\mathbf{h}_t) = \mathbf{y}(\mathbf{x}_{1:t})$$

- network has (causal) memory encoded in **history vector** ( $\mathbf{h}_t$ )





- Variational: introduce **latent** variable sequence  $\mathbf{z}_{1:T}$

$$\begin{aligned}
 p(\mathbf{y}_t | \mathbf{x}_{1:t}) &\approx \int p(\mathbf{y}_t | \mathbf{x}_t, \mathbf{z}_t, \mathbf{h}_{t-1}) p(\mathbf{z}_t | \mathbf{h}_{t-1}) d\mathbf{z}_t \\
 &\approx \int p(\mathbf{y}_t | \mathbf{h}_t) p(\mathbf{z}_t | \mathbf{h}_{t-1}) d\mathbf{z}_t
 \end{aligned}$$

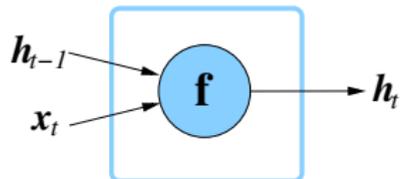
- $\mathbf{z}_t$  a function of complete history (complicates training)

- A flexible extension to activation function is **gating**
  - standard form is ( $\sigma()$  **sigmoid** activation function)

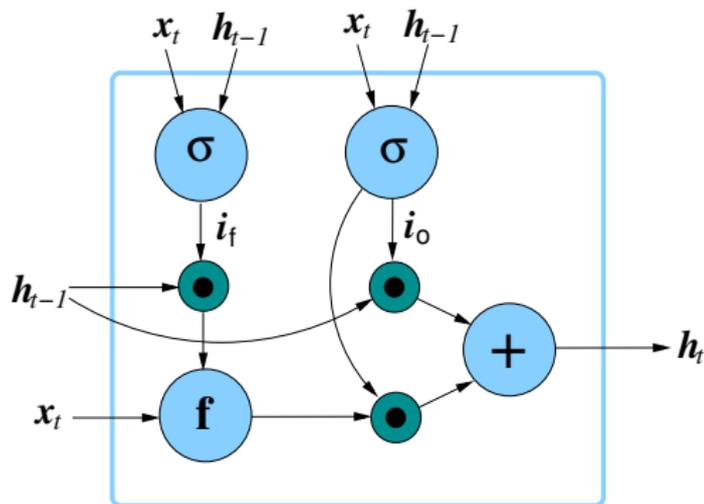
$$\mathbf{i} = \sigma(\mathbf{W}^f \mathbf{x}_t + \mathbf{W}^r \mathbf{h}_{t-1} + \mathbf{b})$$

- vector acts a probabilistic gate on network values
- Gating can be applied at various levels
  - **features**: impact of input/output features on nodes
  - **time**: memory of the network
  - **layer**: influence of a layer's activation function

# Gated Recurrent Unit [6]



Recurrent unit



Gated Recurrent Unit

- Gated Recurrent Unit (GRU) introduces
  - forget gate ( $i_f$ ): gating over time
  - output gate ( $i_o$ ): gating over features (and time)
- Relationships (standard configuration - there are variants)

$$i_f = \sigma(\mathbf{W}_f^f \mathbf{x}_t + \mathbf{W}_f^r \mathbf{h}_{t-1} + \mathbf{b}_f)$$

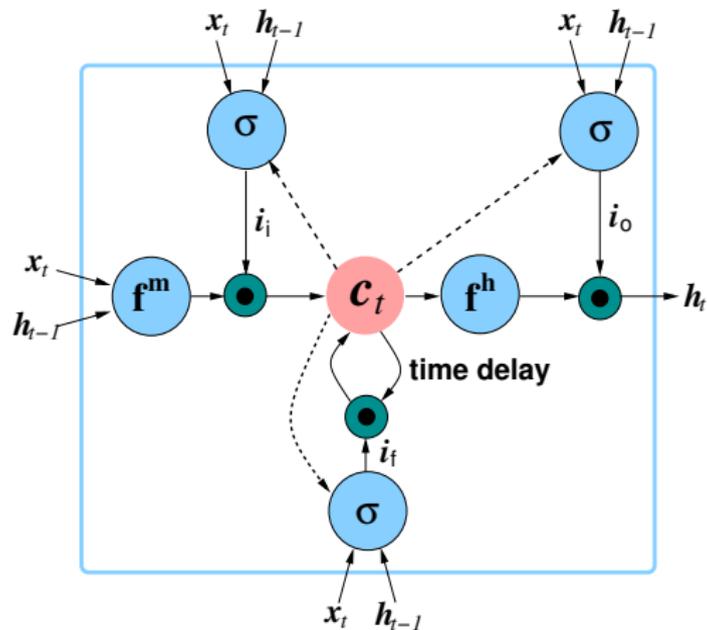
$$i_o = \sigma(\mathbf{W}_o^f \mathbf{x}_t + \mathbf{W}_o^r \mathbf{h}_{t-1} + \mathbf{b}_o)$$

$$\mathbf{y}_t = \mathbf{f}(\mathbf{W}_y^f \mathbf{x}_t + \mathbf{W}_y^r (i_f \odot \mathbf{h}_{t-1}) + \mathbf{b}_y)$$

$$\mathbf{h}_t = i_o \odot \mathbf{h}_{t-1} + (\mathbf{1} - i_o) \odot \mathbf{y}_t$$

- $\odot$  represents **element-wise** multiplication between vectors

# Long-Short Term Memory Networks (reference) [13, 10]



# Long-Short Term Memory Networks (reference)

- The operations can be written as (peephole config):
  - Forget gate ( $i_f$ ), Input gate ( $i_i$ ), Output gate ( $i_o$ )

$$i_f = \sigma(\mathbf{W}_f^f \mathbf{x}_t + \mathbf{W}_f^r \mathbf{h}_{t-1} + \mathbf{W}_f^m \mathbf{c}_{t-1} + \mathbf{b}_f)$$

$$i_i = \sigma(\mathbf{W}_i^f \mathbf{x}_t + \mathbf{W}_i^r \mathbf{h}_{t-1} + \mathbf{W}_i^m \mathbf{c}_{t-1} + \mathbf{b}_i)$$

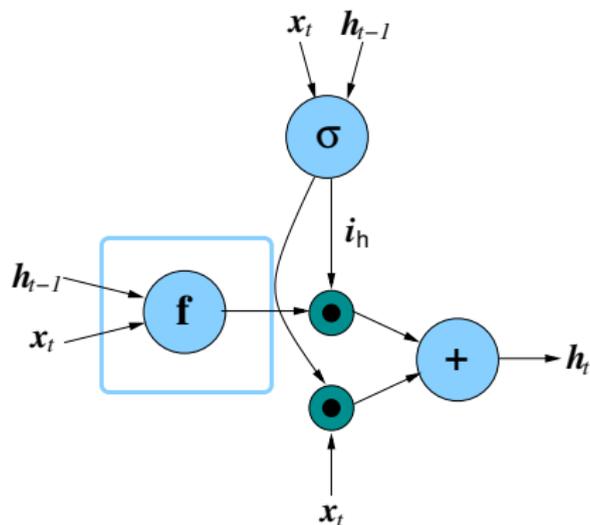
$$i_o = \sigma(\mathbf{W}_o^f \mathbf{x}_t + \mathbf{W}_o^r \mathbf{h}_{t-1} + \mathbf{W}_o^m \mathbf{c}_t + \mathbf{b}_o)$$

- Memory Cell, history vector and gates are related by

$$\mathbf{c}_t = i_f \odot \mathbf{c}_{t-1} + i_i \odot \mathbf{f}^m(\mathbf{W}_c^f \mathbf{x}_t + \mathbf{W}_c^r \mathbf{h}_{t-1} + \mathbf{b}_c)$$

$$\mathbf{h}_t = i_o \odot \mathbf{f}^h(\mathbf{c}_t)$$

- more complicated than GRU (three gates, memory cell)
- memory cell weight matrices ( $\mathbf{W}_f^m, \mathbf{W}_i^m, \mathbf{W}_o^m$ ) diagonal
- can allow explicit analysis of individual cell elements

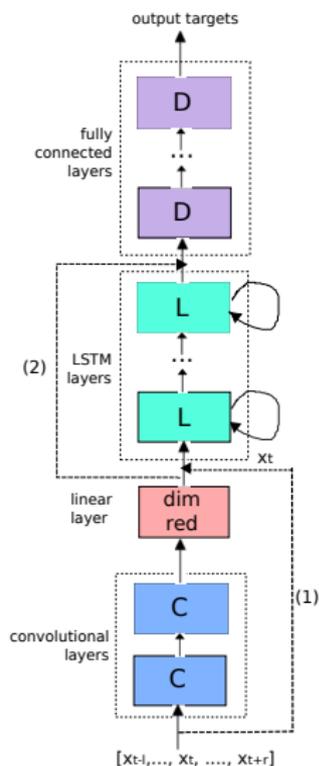


- Gate the output of the node (example from recurrent unit)
  - combine with output from previous layer ( $x_t$ )

$$i_h = \sigma(\mathbf{W}_1^f \mathbf{x}_t + \mathbf{W}_1^r \mathbf{h}_{t-1} + \mathbf{b}_1)$$

$$\mathbf{h}_t = i_h \odot \mathbf{f}(\mathbf{W}_h^f \mathbf{x}_t + \mathbf{W}_h^r \mathbf{h}_{t-1} + \mathbf{b}_h) + (\mathbf{1} - i_h) \odot \mathbf{x}_t$$

# Example Deep Architecture: ASR (reference) [22]



- Example Architecture from Google (2015)
  - C: CNN layer (with pooling)
  - L: LSTM layer
  - D: fully connected layer
- Two multiple layer “skips”
  - (1) connects input to LSTM input
  - (2) connects CNN output to DNN input
- Additional linear projection layer
  - reduces dimensionality
  - and number of network parameters!

# Network Training and Error Back Propagation

# Training Data: Classification

- Supervised training data comprises
  - $\mathbf{x}_i$ :  $d$ -dimensional training observation
  - $y_i$ : class label,  $K$  possible (discrete) classes
- Encode class labels as 1-of- $K$  (“one-hot”) coding:  $y_i \rightarrow \mathbf{t}_i$ 
  - $\mathbf{t}_i$  is the  $K$ -dimensional target vector for  $\mathbf{x}_i$
  - zero other than element associated with class-label  $y_i$
- Consider a network with parameters  $\theta$  and training examples:

$$\{\{\mathbf{x}_1, \mathbf{t}_1\} \dots, \{\mathbf{x}_n, \mathbf{t}_n\}\}$$

- need “distance” from target  $\mathbf{t}_i$  to network output  $\mathbf{y}(\mathbf{x}_i)$

- Least squares error: one of the most common training criteria.

$$E(\boldsymbol{\theta}) = \frac{1}{2} \sum_{p=1}^n \|\mathbf{y}(\mathbf{x}_p) - \mathbf{t}_p\|^2 = \frac{1}{2} \sum_{p=1}^n \sum_{i=1}^K (y_i(\mathbf{x}_p) - t_{pi})^2$$

- Cross-Entropy: note non-zero minimum (entropy of targets)

$$E(\boldsymbol{\theta}) = - \sum_{p=1}^n \sum_{i=1}^K t_{pi} \log(y_i(\mathbf{x}_p))$$

- Cross-Entropy for two classes: single binary target

$$E(\boldsymbol{\theta}) = - \sum_{p=1}^n (t_p \log(y(\mathbf{x}_p)) + (1 - t_p) \log(1 - y(\mathbf{x}_p)))$$

- Supervised training data comprises
  - $\mathbf{x}_i$ :  $d$ -dimensional training observation
  - $\mathbf{y}_i$ :  $K$ -dimensional (continuous) output vector
- Consider a network with parameters  $\theta$  and training examples:

$$\{\{\mathbf{x}_1, \mathbf{y}_1\}, \dots, \{\mathbf{x}_n, \mathbf{y}_n\}\}$$

- need “distance” from target  $\mathbf{y}_i$  to network output  $\mathbf{y}(\mathbf{x}_i)$
- Least squares commonly used criterion

$$E(\theta) = \frac{1}{2} \sum_{p=1}^n (\mathbf{y}(\mathbf{x}_p) - \mathbf{y}_p)' (\mathbf{y}(\mathbf{x}_p) - \mathbf{y}_p)$$

- $\mathbf{y}(\mathbf{x}_i)$  may be viewed as the mean of the prediction

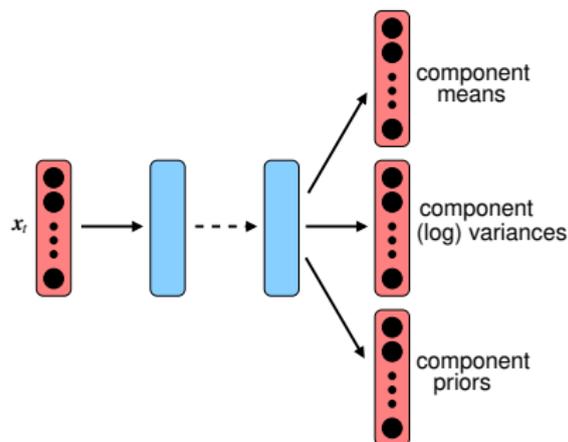
- Generalise least squares (LS) to **maximum likelihood** (ML)

$$E(\boldsymbol{\theta}) = \sum_{p=1}^n \log(p(\mathbf{y}_p | \mathbf{x}_p, \boldsymbol{\theta}))$$

- LS is ML with a single Gaussian, identity covariance matrix
- Criterion appropriate to deep learning for generative models
- Output-layer activation function to ensure valid distribution
  - consider the case of the variance  $\sigma > 0$
  - apply an **exponential** activation function for variances

$$\exp(y_i(\mathbf{x})) > 0$$

- for means just use a **linear** activation function

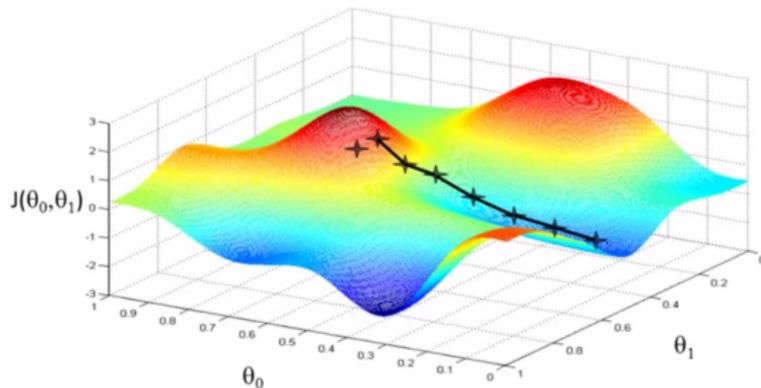


- Predict a **mixture** of  $M$  Gaussians
  - $\mathcal{F}_m^{(c)}(\mathbf{x}_t)$ : prior prediction
  - $\mathcal{F}_m^{(\mu)}(\mathbf{x}_t)$ : mean prediction
  - $\mathcal{F}_m^{(\sigma)}(\mathbf{x}_t)$ : variance prediction
- For component  $m$ , output

$$\mathbf{y}_m(\mathbf{x}_t) = \begin{bmatrix} \mathcal{F}_m^{(c)}(\mathbf{x}_t) \\ \mathcal{F}_m^{(\mu)}(\mathbf{x}_t) \\ \mathcal{F}_m^{(\sigma)}(\mathbf{x}_t) \end{bmatrix}$$

- Optimise using **maximum likelihood** where

$$p(\mathbf{y}_t | \mathbf{x}_t) = \sum_{m=1}^M \mathcal{F}_m^{(c)}(\mathbf{x}_t) \mathcal{N}(\mathbf{y}_t; \mathcal{F}_m^{(\mu)}(\mathbf{x}_t), \mathcal{F}_m^{(\sigma)}(\mathbf{x}_t))$$



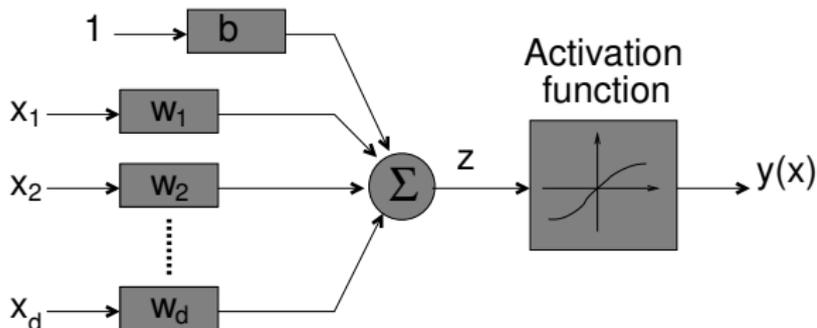
- If there is no closed-form solution - use **gradient descent**

$$\boldsymbol{\theta}[\tau + 1] = \boldsymbol{\theta}[\tau] - \boldsymbol{\Delta}\boldsymbol{\theta}[\tau] = \boldsymbol{\theta}[\tau] - \eta \left. \frac{\partial E}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}[\tau]}$$

- how to get the gradient for all model parameters
- how to avoid **local minima**
- need to consider how to set  $\eta$

- Networks usually have a large number of hidden layers ( $L$ )
  - enables network to model highly non-linear, complex, mappings
  - complicates the training process of the network parameters
- Network parameters are (usually) weights for each layer
  - merging bias vector for each layer into weight matrix

$$\theta = \{ \mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L+1)} \}$$



- Initially just consider a **single layer perceptron**

# Single Layer Perceptron Training

- Take the example of **least squares error** cost function

$$E(\boldsymbol{\theta}) = \frac{1}{2} \sum_{p=1}^n \|y(\mathbf{x}_p) - t_p\|^2 = \sum_{p=1}^n E^{(p)}(\boldsymbol{\theta})$$

- Use chain rule to compute derivatives through network

$$\frac{\partial E(\boldsymbol{\theta})}{\partial w_i} = \left( \frac{\partial E(\boldsymbol{\theta})}{\partial y(\mathbf{x})} \right) \left( \frac{\partial y(\mathbf{x})}{\partial z} \right) \left( \frac{\partial z}{\partial w_i} \right)$$

- change of error function with network output (cost function)
  - change of network output with  $z$  (activation function)
  - change of  $z$  with network weight (parameter to estimate)
- For a **sigmoid** activation function this yields

$$\frac{\partial E^{(p)}(\boldsymbol{\theta})}{\partial w_i} = (y(\mathbf{x}_p) - t_p) y(\mathbf{x}_p) (1 - y(\mathbf{x}_p)) x_{pi}$$

# Error Back Propagation

- Simple concept can be extended to multiple (hidden) layers
  - output from layer  $k - 1$  ( $\mathbf{y}^{(k-1)}$ ) is the input to layer  $k$  ( $\mathbf{x}^{(k)}$ )
- $L + 1$  layer network - use **backward** recursion (see notes)
  - model parameters  $\theta = \{ \mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L+1)} \}$
  - final output  $\mathbf{y}(\mathbf{x}) = \mathbf{y}^{(L+1)}$

$$\delta^{(k)} = \mathbf{\Lambda}^{(k)} \mathbf{W}^{(k+1)'} \delta^{(k+1)}; \quad \frac{\partial E(\theta)}{\partial \mathbf{W}^{(k)}} = \delta^{(k)} \mathbf{x}^{(k)'}$$

- $\mathbf{\Lambda}^{(k)} = \frac{\partial \mathbf{y}^{(k)}}{\partial \mathbf{z}^{(k)}}$ : layer  $k$  activation derivative matrix
- $\mathbf{W}^{(k+1)'}$ : weight matrix for layer  $k + 1$
- $\delta^{(k+1)}$  =  $\frac{\partial E(\theta)}{\partial \mathbf{z}^{(k+)}}$ : error vector for layer  $k + 1$

# Error Back Propagation (cont)

- The process to get the derivative involves:
  1. For input vector  $\mathbf{x}_p$  propagate **forward**: yields ( $1 \leq k \leq L$ )
    - $\mathbf{y}^{(k)}$  the output value for each node of layer all layers
    - $\mathbf{z}^{(k)}$  the input value to the non-linearity for layer  $k$
  2. Compute  $\left. \frac{\partial E(\boldsymbol{\theta})}{\partial \mathbf{y}(\mathbf{x})} \right|_{\boldsymbol{\theta}_{[\tau]}}$  (the gradient at the output layer).
  3. From the output gradient propagate **backwards**: yields
    - $\boldsymbol{\delta}^{(k)}$  the error vector for each layer
    - $\frac{\partial E(\boldsymbol{\theta})}{\partial \mathbf{W}^{(k)}}$ : the (desired) derivative for layer  $k$  weights

# Optimisation

# Batch/On-Line Gradient Descent

- The default gradient is computed over all samples
  - for large data sets very slow - each update
- Modify to batch update - just use a subset of data,  $\tilde{\mathcal{D}}$ ,

$$E(\theta) = - \sum_{p \in \tilde{\mathcal{D}}} \sum_{i=1}^K t_{pi} \log(y_i(\mathbf{x}_p))$$

- How to select the subset,  $\tilde{\mathcal{D}}$ ?
  - small subset “poor” estimate of true gradient
  - large subset each parameter update is expensive
- One extreme is to update after each sample
  - $\tilde{\mathcal{D}}$  comprises a single sample in order
  - “noisy” gradient estimate for updates

# Stochastic Gradient Descent (SGD)

- Two modifications to the baseline approaches
  1. Randomise the order of the data presented for training
    - important for structured data
  2. Introduce **mini-batch** updates
    - $\tilde{\mathcal{D}}$  is a (random) subset of the training data
    - better estimate of the gradient at each update
    - **but** reduces number of iterations
- Mini-batch updates are (almost) always used
  - make use of parallel processing (GPUs) for efficiency
- Research of parallel versions of SGD on-going

- A number of issues for gradient descent including:
  - stops at local maxima
  - handling “ravines”
- **Momentum** aims to address this - parameter change becomes:

$$\Delta\theta[\tau] = \eta \left. \frac{\partial E(\theta)}{\partial \theta} \right|_{\theta[\tau]} + \alpha \Delta\theta[\tau - 1]$$

- smooths parameter changes over iterations
  - introduces an additional tunable parameter
- For simplicity introduce compact notation

$$\left. \frac{\partial E(\theta)}{\partial \theta} \right|_{\theta[\tau]} = \nabla(E(\theta[\tau]))$$

# Adaptive Learning Rates

- Speed of convergence depends on  $\eta$ 
  - **too large**: updates may diverge rather than converge
  - **too small**: very slow convergence (impractical)
- The standard expression has a fixed learning rate
  - can we have learning rate change with iteration

$$\boldsymbol{\theta}[\tau + 1] = \boldsymbol{\theta}[\tau] - \Delta\boldsymbol{\theta}[\tau] = \boldsymbol{\theta}[\tau] - \eta[\tau]\nabla(E(\boldsymbol{\theta}[\tau]))$$

- how to set  $\eta[\tau]$  (or generally parameter update  $\Delta\boldsymbol{\theta}[\tau]$ )?
- One very simple approach

$$\eta[\tau + 1] = \begin{cases} 1.1\eta[\tau]; & \text{if } E(\boldsymbol{\theta}[\tau]) < E(\boldsymbol{\theta}[\tau - 1]) \\ 0.5\eta[\tau]; & \text{if } E(\boldsymbol{\theta}[\tau]) > E(\boldsymbol{\theta}[\tau - 1]) \end{cases}$$

- increase learning rate when going in “correct direction”

## Gradient Descent Refinements (reference)

- **Nesterov**: concept of gradient at next iteration

$$\Delta\theta[\tau] = \eta \nabla(E(\theta[\tau] - \alpha \Delta\theta[\tau - 1])) + \alpha \Delta\theta[\tau - 1]$$

- **AdaGrad**: dimension specific learning rates ( $\epsilon$  floor parameter)

$$\Delta\theta[\tau] = \eta \beta_t \odot \left. \frac{\partial E(\theta)}{\partial \theta} \right|_{\theta[\tau]} ; \quad \beta_{ti} = \frac{1}{\sqrt{\epsilon + \sum_{t=1}^{\tau} \nabla_i(E(\theta[t]))^2}}$$

- $\epsilon$  is a smoothing term to avoid division by zero
- **Adam**: Adaptive Moment Estimation: use dimension moments

$$\Delta\theta_i[\tau] = \frac{\eta}{\sqrt{\sigma_{\tau i}^2 + \epsilon}} \mu_{\tau i}; \quad \begin{aligned} \mu_{\tau i} &= \alpha_1 \mu_{(\tau-1)i} + (1 - \alpha_1) \nabla_i(E(\theta[\tau])) \\ \sigma_{\tau i}^2 &= \alpha_2 \sigma_{(\tau-1)i}^2 + (1 - \alpha_2) \nabla_i(E(\theta[\tau]))^2 \end{aligned}$$

- additional normalisation applied to  $\mu_{\tau i}$  and  $\sigma_{\tau i}^2$  to offset initialisation bias

## Second-Order Approximations

- Gradient descent makes use of first-order derivatives of
  - what about **higher order derivatives**? Consider

$$E(\boldsymbol{\theta}) = E(\boldsymbol{\theta}[\tau]) + (\boldsymbol{\theta} - \boldsymbol{\theta}[\tau])' \mathbf{g} + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}[\tau])' \mathbf{H} (\boldsymbol{\theta} - \boldsymbol{\theta}[\tau]) + \mathcal{O}(\boldsymbol{\theta}^3)$$

where

$$\mathbf{g} = \nabla E(\boldsymbol{\theta}[\tau]); \quad (\mathbf{H})_{ij} = h_{ij} = \left. \frac{\partial^2 E(\boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} \right|_{\boldsymbol{\theta}[\tau]}$$

- Ignoring higher order terms and equating to zero

$$\nabla E(\boldsymbol{\theta}) = \mathbf{g} + \mathbf{H}(\boldsymbol{\theta} - \boldsymbol{\theta}[\tau])$$

Equating to zero (check minimum!) -  $\mathbf{H}^{-1} \mathbf{g}$  **Newton direction**

$$\boldsymbol{\theta}[\tau + 1] = \boldsymbol{\theta}[\tau] - \mathbf{H}^{-1} \mathbf{g}; \quad \Delta \boldsymbol{\theta}[\tau] = \mathbf{H}^{-1} \mathbf{g}$$

## Issues with Second-Order Approaches

1. The evaluation of the Hessian may be computationally expensive as  $\mathcal{O}(N^2)$  parameters must be accumulated for each of the  $n$  training samples.
2. The Hessian must be inverted to find the direction,  $\mathcal{O}(N^3)$ . This gets very expensive as  $N$  gets large.
3. The direction given need not head towards a minimum - it could head towards a **maximum** or **saddle point**. This occurs if the Hessian is not **positive-definite** i.e.

$$\mathbf{v}'\mathbf{H}\mathbf{v} > 0$$

for all  $\mathbf{v}$ . The Hessian may be made positive definite using

$$\tilde{\mathbf{H}} = \mathbf{H} + \lambda\mathbf{I}$$

If  $\lambda$  is large enough then  $\tilde{\mathbf{H}}$  is positive definite.

4. If the surface is highly non-quadratic the step sizes may be too large and the optimisation becomes unstable.

- Interesting making use of the error curvature, assumptions:
  - error surface is quadratic in nature
  - weight gradients treated independently (diagonal Hessian)
- Using these assumptions

$$E(\theta) \approx E(\theta[\tau]) + b(\theta - \theta[\tau]) + a(\theta - \theta[\tau])^2$$
$$\frac{\partial E(\theta)}{\partial \theta} \approx b + 2a(\theta - \theta[\tau])$$

- To find  $a$  and  $b$  make use of:
  - update step,  $\Delta\theta[\tau - 1]$ , and gradient,  $g[\tau - 1]$ , iteration  $\tau - 1$
  - the gradient at iteration  $\tau$  is  $g[\tau]$
  - after new update  $\Delta\theta[\tau]$  the gradient should be zero
- The following equalities are obtained

$$g[\tau - 1] = b - 2a\Delta\theta[\tau - 1], \quad 0 = b + 2a\Delta\theta[\tau], \quad g[\tau] = b$$

$$\rightarrow \Delta\theta[\tau] = \frac{g[\tau]}{g[\tau - 1] - g[\tau]} \Delta\theta[\tau - 1]$$

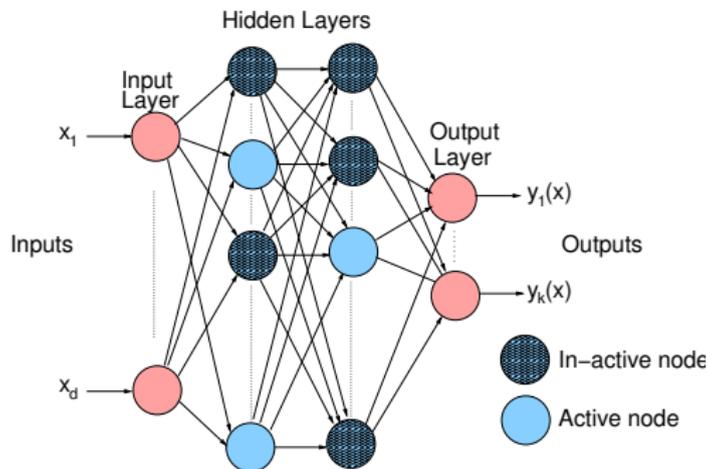


- A major issues with training networks is **generalisation**
- Simplest approach is **early stopping**
  - don't wait for convergence - just stop ...
- To address this forms of **regularisation** are used
  - one standard form is ( $N$  is the total number of weights):

$$\tilde{E}(\theta) = E(\theta) + \nu\Omega(\theta); \quad \Omega(\theta) = \frac{1}{2} \sum_{i=1}^N w_i^2$$

- a zero "prior" is used for the model parameters
- Simple to include in gradient-descent optimisation

$$\nabla \tilde{E}(\theta[\tau]) = \nabla E(\theta[\tau]) + \nu \mathbf{w}[\tau]$$



- Dropout is simple way of improving generalisation
  - randomly de-activate (say) 50% of the nodes in the network
  - update the model parameters
- Prevents a single node specialising to a task

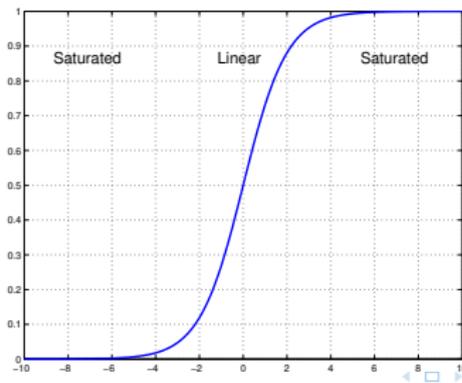
- As with standard classifiers two stage classification often used
  - features are designed by expert
  - current trend to remove two-stage process **end-to-end**
- Features may have different **dynamic ranges**
  - consider dimension 1:  $-1000 \rightarrow 1000$  vs dimension 2:  $-1 \rightarrow 1$
  - can influence “importance” of features at start of training
- Data **whitening** often employed

$$\tilde{x}_{pi} = \frac{x_{pi} - \mu_i}{\sigma_i}; \quad \mu_i = \frac{1}{n} \sum_{p=1}^n x_{pi} \quad \sigma_i^2 = \frac{1}{n} \sum_{p=1}^n (x_{pi} - \mu_i)^2$$

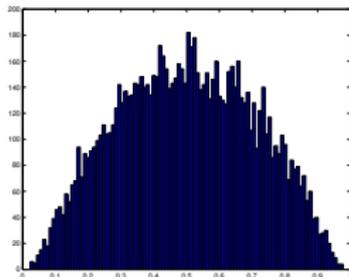
- only influences initialisation (linear transform and bias)

# Network Initialisation: Weight Parameters

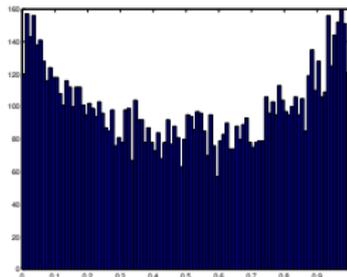
- A starting point (initialisation) for gradient descent is useful
  - one of the “old” concerns with deep networks was initialisation
  - recurrent neural networks are **very** deep!
- It is not possible to guarantee a good starting point, **but**
  - would like a **parsimonious** initialisation
- What about **Gaussian** random initialisation
  - consider zero mean distribution, scale the variance
  - sigmoid non-linearity



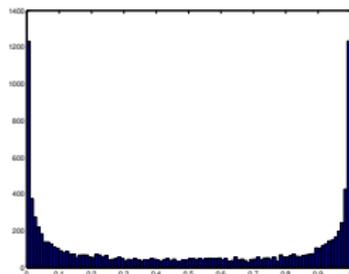
# Gaussian Initialisation



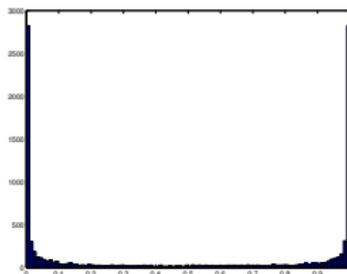
$$\mathcal{N}(0, 1)$$



$$\mathcal{N}(0, 2)$$



$$\mathcal{N}(0, 4)$$



$$\mathcal{N}(0, 8)$$

- Pass 1-dimensional data through a [sigmoid](#)

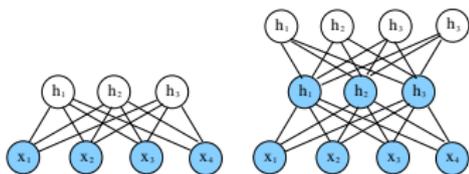
- Need to worry about the following gradient issues
  - **Vanishing**: derivatives go to zero - parameters not updated
  - **Exploding**: derivatives get very large - cause saturation
- **Xavier Initialisation**: simple scheme for initialising weights
  - linear activation functions  $\mathbf{y} = \mathbf{W}\mathbf{x}$ ,
  - assuming all weights/observations independent
  - $\mathbf{x}$   $n$  dimensional, zero mean, identity variance

$$\text{Var}(y_i) = \text{Var}(\mathbf{w}'_i \mathbf{x}) = n \text{Var}(w_{ij}) \text{Var}(x_i)$$

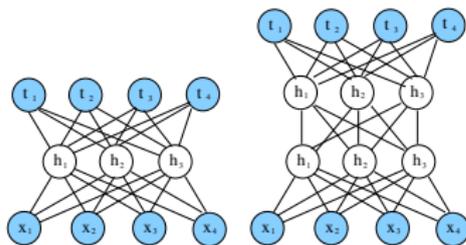
- Would like variance on output to be the same as the input

$$\text{Var}(w_{ij}) = \frac{1}{n} \quad (\text{for : } \text{Var}(y_i) = \text{Var}(x_i))$$

# Alternative Initialisation schemes



Generative Pre-Training



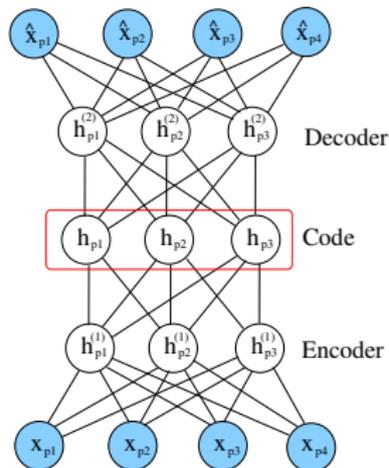
Discriminative Pre-Training

- Restricted Boltzmann machine: generative pre-training
  - initialise the network parameters using a generative RBM
  - train the RBM layer-by-layer
  - if only the output layer trained this is a **deep belief network**
- Layer-by-layer training: discriminative pre-training
  - iterative process:
    1. remove o/p layer, add a random initialised layer, add o/p layer
    2. train (limited iterations e.g. 1) network and then goto (1)

# Example Systems (reference)

# Autoencoders (Non-Linear Feature Extraction)

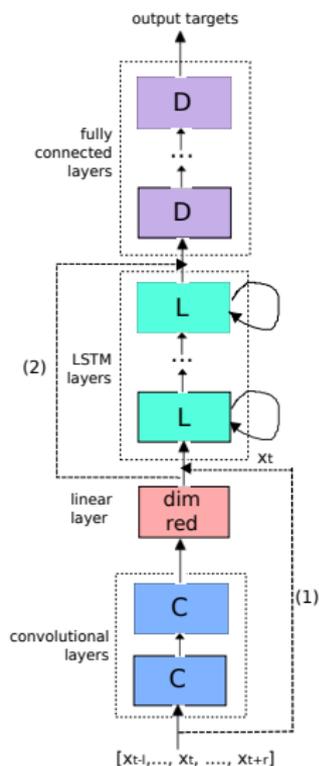
- An autoencoder is a particular form of feed-forward network
  - a (often low-dimensional) **code** layer ( $\mathbf{h}_p$ )
  - trained to reproduce the input at the output



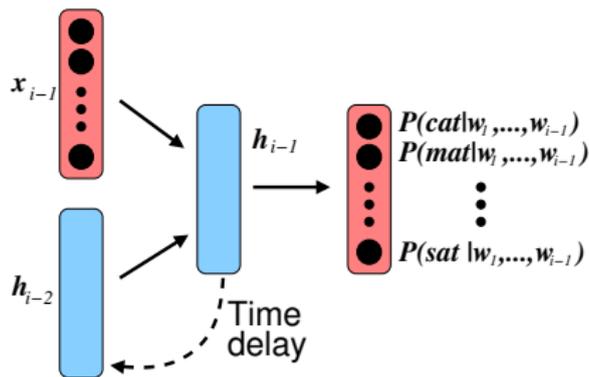
- Training criterion

$$E(\boldsymbol{\theta}) = \sum_{p=1}^n f(\mathbf{x}_p, \hat{\mathbf{x}}_p)$$

- Can be used to **denoise** data
  - “noise-corrupted” data in,
  - distance to “clean” data for training



- Example Architecture from Google (2015)
  - C: CNN layer (with pooling)
  - L: LSTM layer
  - D: fully connected layer
- Two multiple layer “skips”
  - (1) connects input to LSTM input
  - (2) connects CNN output to DNN input
- Additional linear projection layer
  - reduces dimensionality
  - and number of network parameters!



- Neural networks extensively used for language modelling
  - compute  $P(\omega_{1:L})$ : used in translation/ASR/topic spotting etc
- 1-of-K (“one-hot”) coding for  $i^{th}$  word,  $\omega_i$ ,  $\mathbf{x}_i$ 
  - additional **out-of-shortlist** symbol may be added
  - **softmax** activation function on output layer

- Neural networks extensively used for language modelling
  - recurrent neural networks - complete word history

$$P(\omega_{1:L}) = \prod_{i=1}^L P(\omega_i | \omega_{1:i-1}) \approx \prod_{i=1}^L P(\omega_i | \omega_{i-1}, \tilde{\mathbf{h}}_{i-2}) \approx \prod_{i=1}^L P(\omega_i | \tilde{\mathbf{h}}_{i-1})$$

- Input and output layer sizes can be very large
  - size of vocabulary ( $> 10K$ ), not an issue for input
  - output-layer (softmax) expensive (normalisation term)
- Issues that need to be addressed
  1. training: how to efficiently train on billions of words?
  2. decoding for ASR: how to handle dependence on complete history?

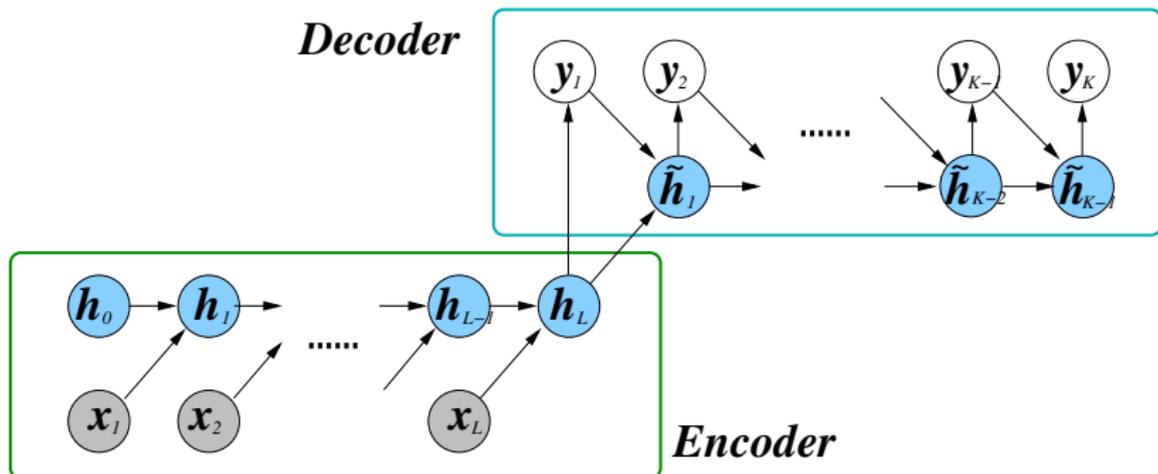
- Train a discriminative model from
  - $\mathbf{x}_{1:L} = \{\mathbf{x}_1, \dots, \mathbf{x}_L\}$ :  $L$ -length input sequence (source language)
  - $\mathbf{y}_{1:K} = \{\mathbf{y}_1, \dots, \mathbf{y}_K\}$ :  $K$ -length output (target language)

$$\begin{aligned} p(\mathbf{y}_{1:K} | \mathbf{x}_{1:L}) &= \prod_{i=1}^K p(\mathbf{y}_i | \mathbf{y}_{1:i-1}, \mathbf{x}_{1:L}) \\ &\approx \prod_{i=1}^L p(\mathbf{y}_i | \mathbf{y}_{i-1}, \tilde{\mathbf{h}}_{i-2}, \mathbf{c}) \end{aligned}$$

- need to map  $\mathbf{x}_{1:L}$  to a fixed-length vector

$$\mathbf{c} = \phi(\mathbf{x}_{1:L})$$

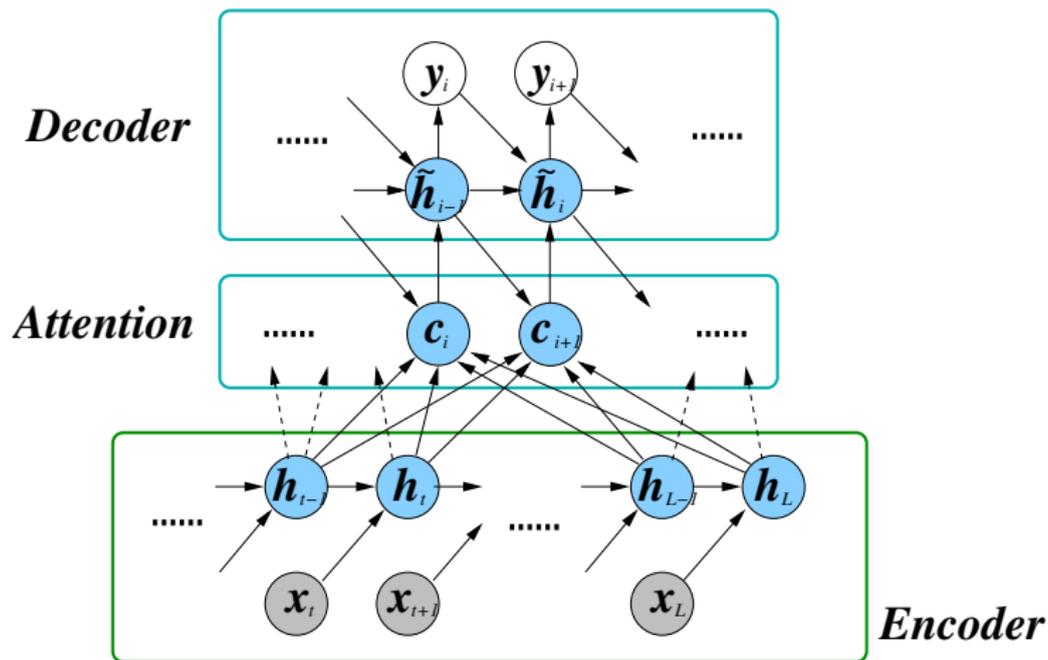
- $\mathbf{c}$  is a fixed length vector - like a [sequence kernel](#)



- One form is to use **hidden unit** from acoustic RNN/LSTM

$$\mathbf{c} = \phi(\mathbf{x}_{1:L}) = \mathbf{h}_L$$

- dependence on context is global via  $\mathbf{c}$  - possibly limiting



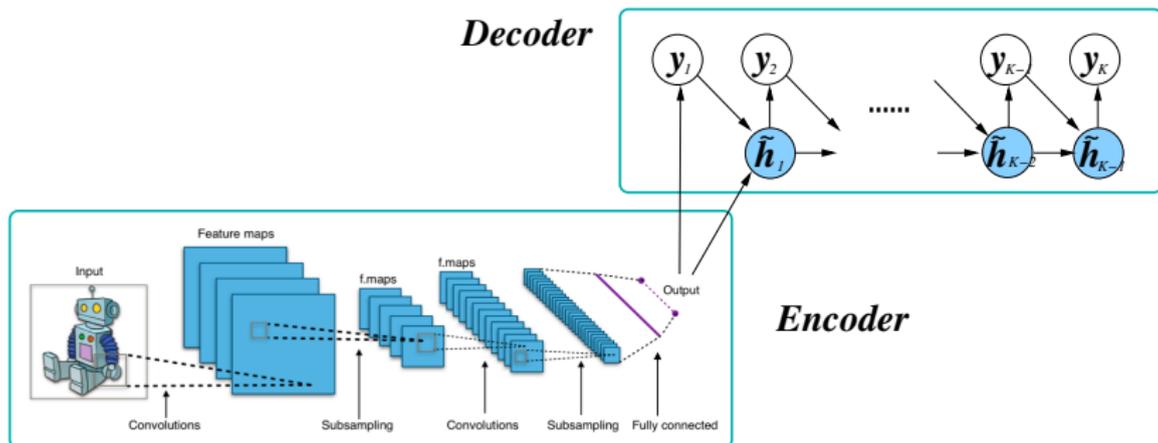
- Introduce **attention** layer to system
  - introduce dependence on locality  $i$

$$p(\mathbf{y}_{1:K}|\mathbf{x}_{1:L}) \approx \prod_{i=1}^K p(\mathbf{y}_i|\mathbf{y}_{i-1}, \tilde{\mathbf{h}}_{i-1}, \mathbf{c}_i) \approx \prod_{i=1}^K p(\mathbf{y}_i|\tilde{\mathbf{h}}_{i-1})$$

$$\mathbf{c}_i = \sum_{\tau=1}^L \alpha_{i\tau} \mathbf{h}_\tau; \quad \alpha_{i\tau} = \frac{\exp(e_{i\tau})}{\sum_{j=1}^L \exp(e_{ij})}, \quad e_{i\tau} = f^e(\tilde{\mathbf{h}}_{i-2}, \mathbf{h}_\tau)$$

- $e_{i\tau}$  how well position  $i-1$  in input matches position  $\tau$  in output
- $\mathbf{h}_\tau$  is representation (RNN) for the input at position  $\tau$

# Image Captioning



- Encode-image as vector use a deep convolutional network
  - generate caption using recurrent network (RNN/LSTM)
  - all parameters optimised (using example image captions)

# Conclusions

# Is Deep Learning the Solution?

- Deep Learning: state-of-the-art performance in range of tasks
  - machine translation, image recognition/captioning,
  - speech recognition/synthesis ...
- Traditionally use **two-stage approach** to build classifier:
  1. **feature extraction**: convert waveform to parametric form
  2. **modelling**: given parameters train model
- Limitations in feature extraction cannot be overcome ...
  - integrate feature extraction into process
  - attempt to directly model/synthesise waveform (WaveNet)
- **BUT**
  - require large quantities of data (research direction)
  - networks are difficult to optimise - tuning required
  - hard to interpret networks to get insights
  - sometimes difficult to learn from previous tasks ...

- [1] O. Abdel-Hamid, A.-R. Mohamed, H. Jiang, and G. Penn, "Applying convolutional neural networks concepts to hybrid NN-HMM model for speech recognition," in *2012 IEEE international conference on Acoustics, speech and signal processing (ICASSP)*. IEEE, 2012, pp. 4277–4280.
- [2] Y. Bengio, R. Ducharme, P. Vincent, and C. Jauvin, "A neural probabilistic language model," *journal of machine learning research*, vol. 3, no. Feb, pp. 1137–1155, 2003.
- [3] C. Bishop, "Mixture density networks," in *Tech. Rep. NCRG/94/004, Neural Computing Research Group, Aston University*, 1994.
- [4] W. Chan, N. Jaitly, Q. V. Le, and O. Vinyals, "Listen, attend and spell," *CoRR*, vol. abs/1508.01211, 2015. [Online]. Available: <http://arxiv.org/abs/1508.01211>
- [5] J. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio, "Attention-based models for speech recognition," *CoRR*, vol. abs/1506.07503, 2015. [Online]. Available: <http://arxiv.org/abs/1506.07503>
- [6] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," *arXiv preprint arXiv:1412.3555*, 2014.
- [7] J. Chung, K. Kastner, L. Dinh, K. Goel, A. C. Courville, and Y. Bengio, "A recurrent latent variable model for sequential data," *CoRR*, vol. abs/1506.02216, 2015. [Online]. Available: <http://arxiv.org/abs/1506.02216>
- [8] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks." in *Aistats*, vol. 9, 2010, pp. 249–256.
- [9] A. Graves, "Sequence transduction with recurrent neural networks," *CoRR*, vol. abs/1211.3711, 2012. [Online]. Available: <http://arxiv.org/abs/1211.3711>
- [10] A. Graves, A.-R. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in *2013 IEEE international conference on acoustics, speech and signal processing*. IEEE, 2013, pp. 6645–6649.
- [11] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *CoRR*, vol. abs/1512.03385, 2015. [Online]. Available: <http://arxiv.org/abs/1512.03385>
- [12] G. Hinton, L. Deng, D. Yu, G. E. Dahl, A.-R. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N. Sainath, et al., "Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups," *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 82–97, 2012.

- [13] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997.
- [14] Y. LeCun and Y. Bengio, "Convolutional networks for images, speech, and time series," *The handbook of brain theory and neural networks*, vol. 3361, no. 10, p. 1995, 1995.
- [15] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, "Backpropagation applied to handwritten zip code recognition," *Neural computation*, vol. 1, no. 4, pp. 541–551, 1989.
- [16] L. Lu, X. Zhang, K. Cho, and S. Renals, "A study of the recurrent neural network encoder-decoder for large vocabulary speech recognition," in *Proc. INTERSPEECH*, 2015.
- [17] T. Mikolov, M. Karafát, L. Burget, J. Cernocký, and S. Khudanpur, "Recurrent neural network based language model." in *Interspeech*, vol. 2, 2010, p. 3.
- [18] V. Nair and G. E. Hinton, "Rectified linear units improve restricted Boltzmann machines," in *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, 2010, pp. 807–814.
- [19] T. Robinson and F. Fallside, "A recurrent error propagation network speech recognition system," *Computer Speech & Language*, vol. 5, no. 3, pp. 259–274, 1991.
- [20] S. Ruder, "An overview of gradient descent optimization algorithms," <http://sebastianruder.com/optimizing-gradient-descent/index.html>, accessed: 2016-10-14.
- [21] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Parallel distributed processing: Explorations in the microstructure of cognition, vol. 1," D. E. Rumelhart, J. L. McClelland, and C. PDP Research Group, Eds. Cambridge, MA, USA: MIT Press, 1986, ch. Learning Internal Representations by Error Propagation, pp. 318–362.
- [22] T. N. Sainath, O. Vinyals, A. Senior, and H. Sak, "Convolutional, long short-term memory, fully connected deep neural networks," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2015, pp. 4580–4584.
- [23] M. Schuster and K. K. Paliwal, "Bidirectional recurrent neural networks," *IEEE Transactions on Signal Processing*, vol. 45, no. 11, pp. 2673–2681, 1997.

- [24] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting." *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [25] R. K. Srivastava, K. Greff, and J. Schmidhuber, "Highway networks," *CoRR*, vol. abs/1505.00387, 2015. [Online]. Available: <http://arxiv.org/abs/1505.00387>
- [26] P. Swietojanski and S. Renals, "Differentiable pooling for unsupervised acoustic model adaptation," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. PP, no. 99, pp. 1–1, 2016.
- [27] C. Wu, P. Karanasou, M. Gales, and K. C. Sim, "Stimulated deep neural network for speech recognition," in *Proceedings interspeech*, 2016.
- [28] H. Zen and A. Senior, "Deep mixture density networks for acoustic modeling in statistical parametric speech synthesis," in *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2014, pp. 3844–3848.
- [29] C. Zhang and P. Woodland, "DNN speaker adaptation using parameterised sigmoid and ReLU hidden activation functions," in *Proc. ICASSP'16*, 2016.