

# M9. (Artificial) Neural Networks (ANNs)

Balaraman Ravindran

PRML Aug-Dec 2021 (BR section)

# Acknowledgment of Sources

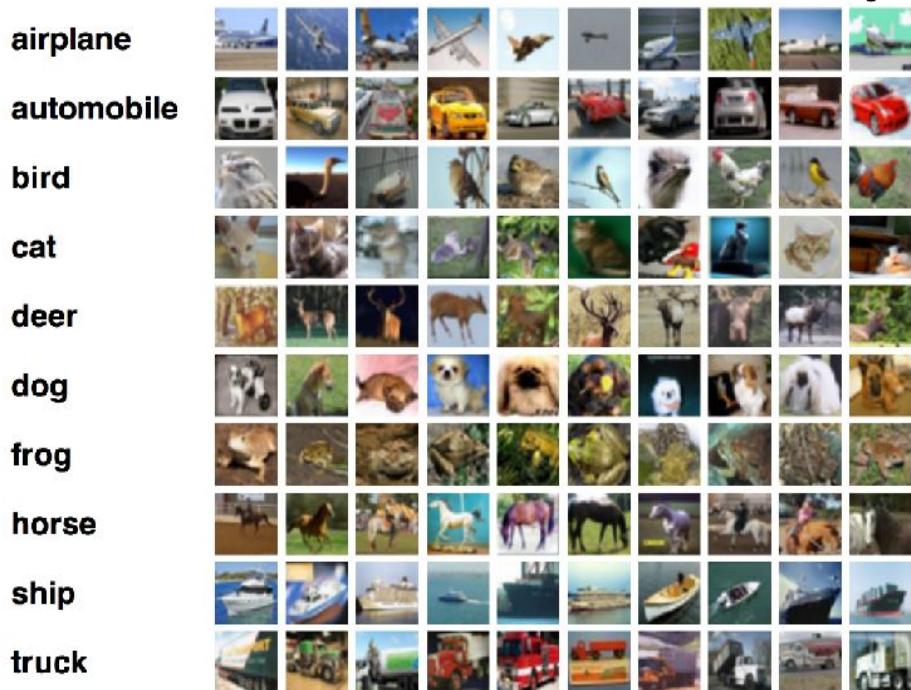
- Slides based on content from related
  - Courses:
    - IITM – Profs. Arun/Harish/Chandra’s PRML offerings (slides, quizzes, notes, etc.), Prof. Ravi’s “Intro to ML” slides – cited respectively as [AR], [HR], [CC], [BR] in the bottom right of a slide.
    - India – NPTEL PR course by IISc Prof. PS. Sastry (slides, etc.) – cited as [PSS] in the bottom right of a slide.
  - Books:
    - PRML by Bishop. (content, figures, slides, etc.) – cited as [CMB]
    - Pattern Classification by Duda, Hart and Stork. (content, figures, etc.) – [DHS]
    - Mathematics for ML by Deisenroth, Faisal and Ong. (content, figures, etc.) – [DFO]
    - Information Theory, Inference and Learning Algorithms by David JC MacKay – [DJM]

# Outline for Module M9

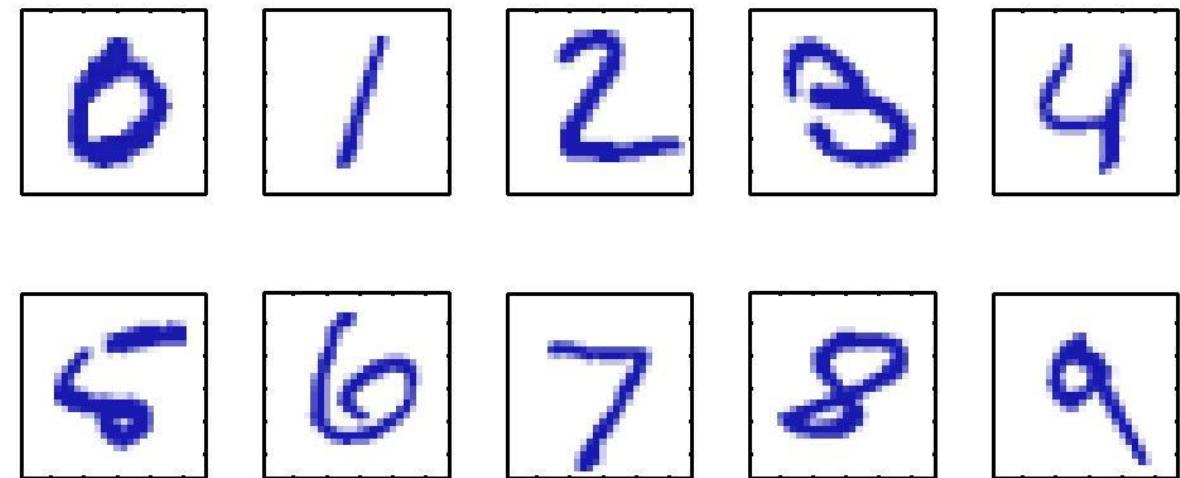
- M9. Neural Networks
  - **M9.0 Introduction/Motivation**
  - M9.1 Feed-forward neural networks
    - (Key Idea: Adaptive Basis functions and its Generalizations)
    - (Network architecture, Network training & Backpropagation algo. sketch)
  - M9.2 Mixture Density Networks
  - M9.3 Concluding thoughts

# Recall: popular examples

Classify an image into one of 10 classes,  
given a “training set” of images with classes.

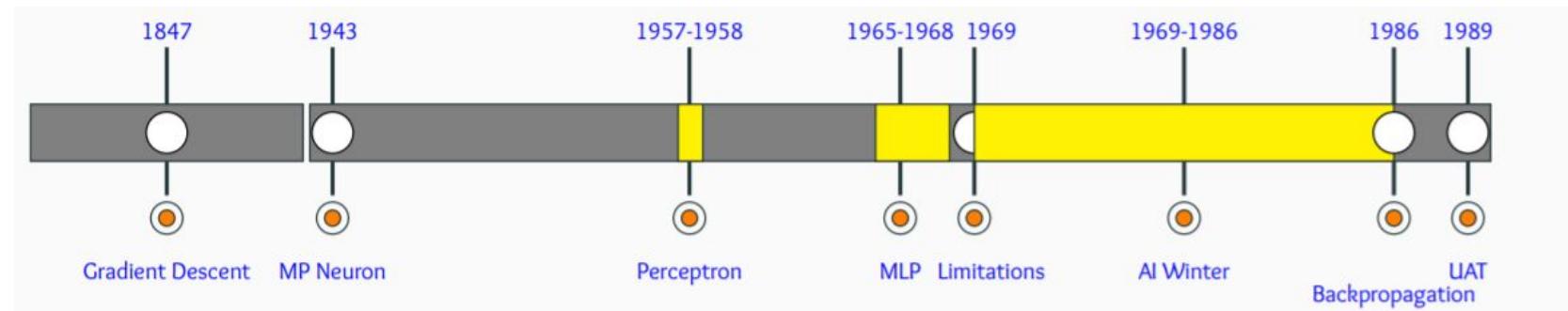
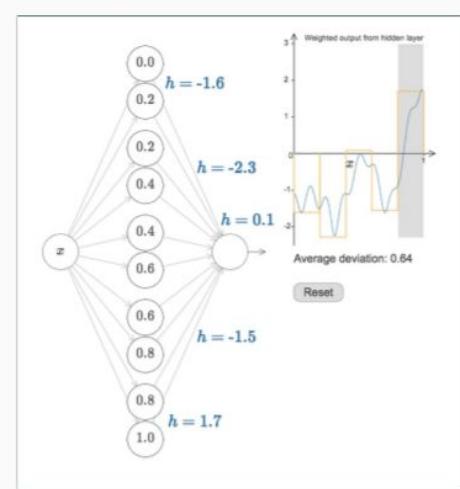
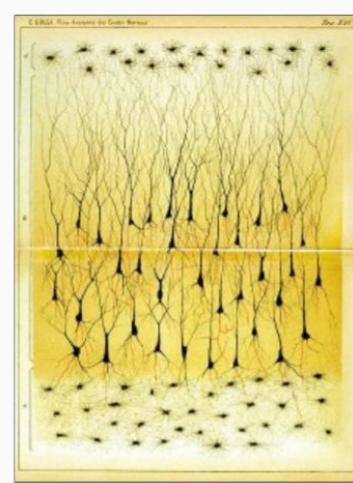
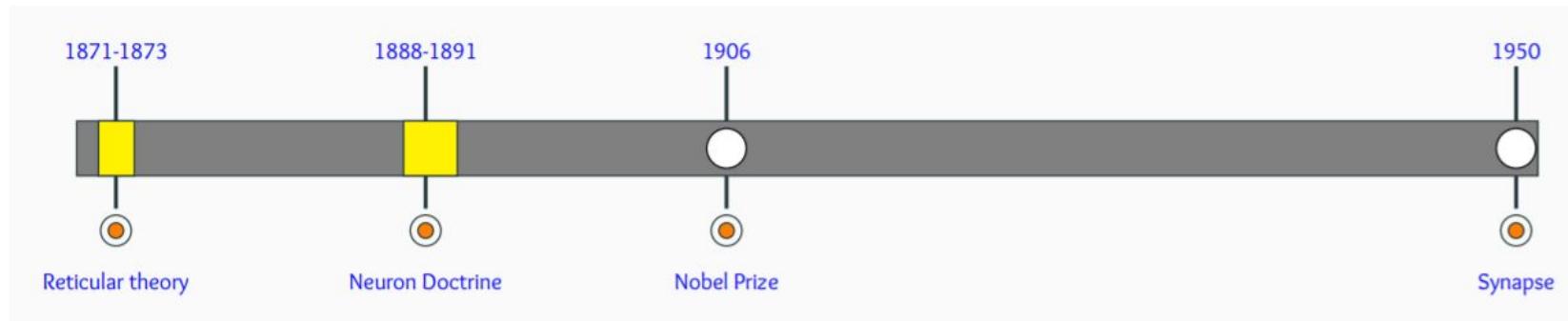


Handwritten Digit Recognition

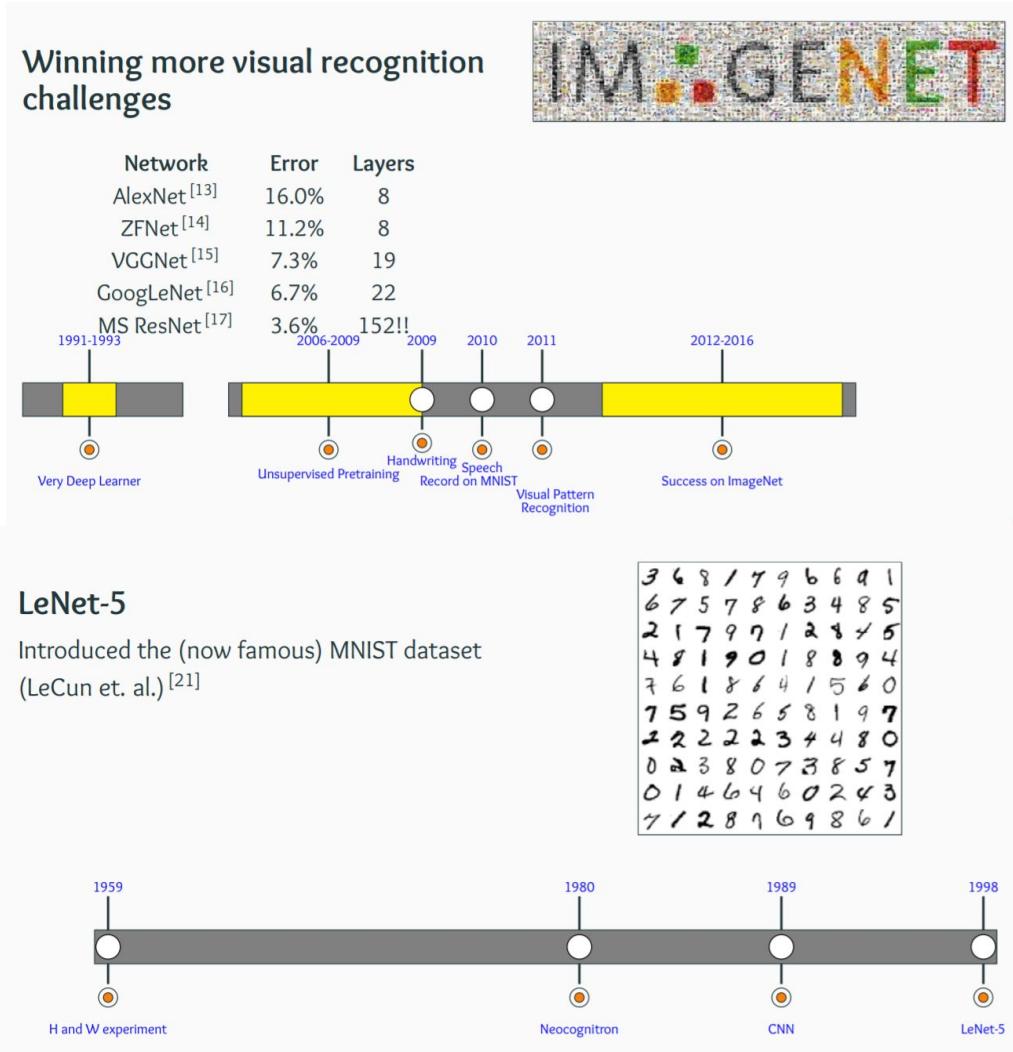


[HR]

# Very brief history of early research on Biological & Artificial Neural Networks



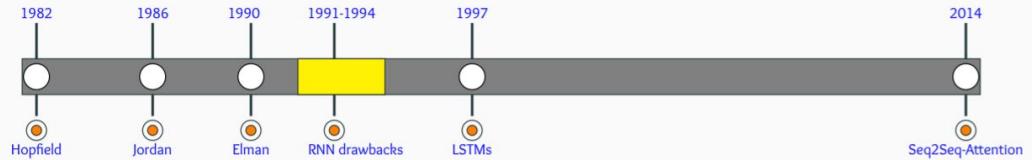
# (Recent) Success stories of ANNs/DL



## Sequence To Sequence Models

Initial success in using RNNs/LSTMs for  
large scale Sequence To Sequence  
Learning Problems

Introduction of Attention which is perhaps the idea of the decade!

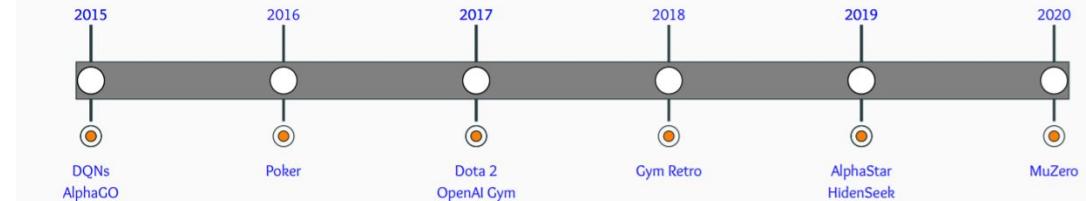


## Jack of all, Master of all!

MuZero masters Go, chess, shogi and Atari without needing to be told the rules, thanks to its ability to plan winning strategies in unknown environments.



<https://deepmind.com/blog>



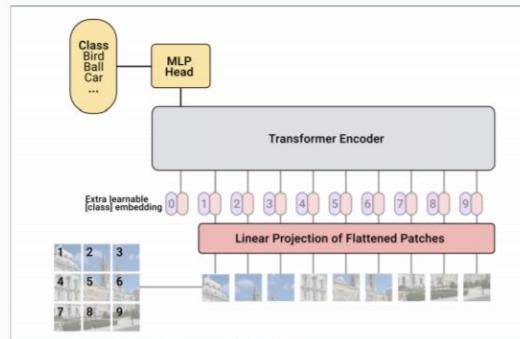
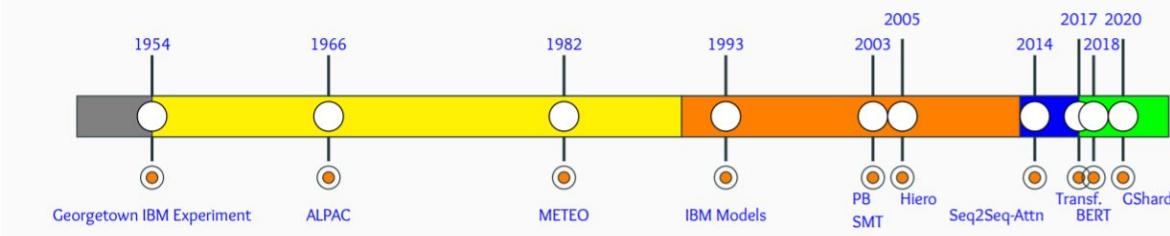
[MK]

# More (recent) success stories of ANNs/DL

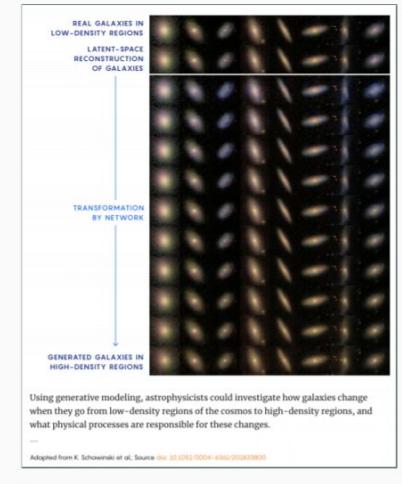
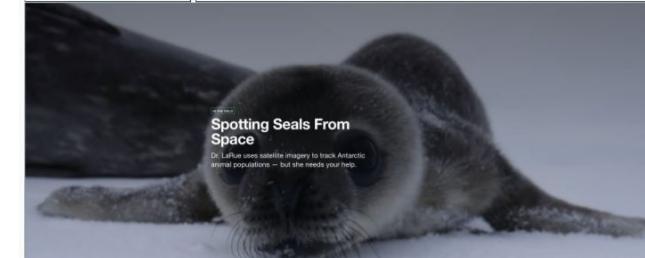
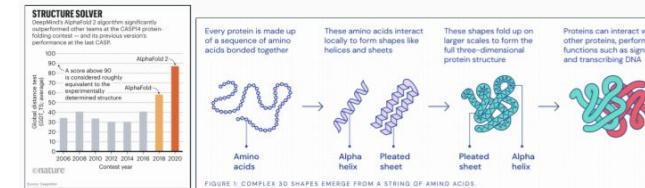
## From Language To Vision

A vision model<sup>a</sup> based as closely as possible on the Transformer architecture originally designed for text-based tasks (another paradigm shift from CNNs which have been around since 1980s!)

<sup>a</sup>[Source:https://ai.googleblog.com/2020/12/transformers-for-image-recognition-at.html](https://ai.googleblog.com/2020/12/transformers-for-image-recognition-at.html)



## Accelerating Scientific Discovery<sup>a</sup>



<sup>a</sup><https://deepmind.com/blog/article/AlphaFold-Using-AI-for-scientific-discovery>

<https://ocean.org/stories/spotting-seals-from-space>

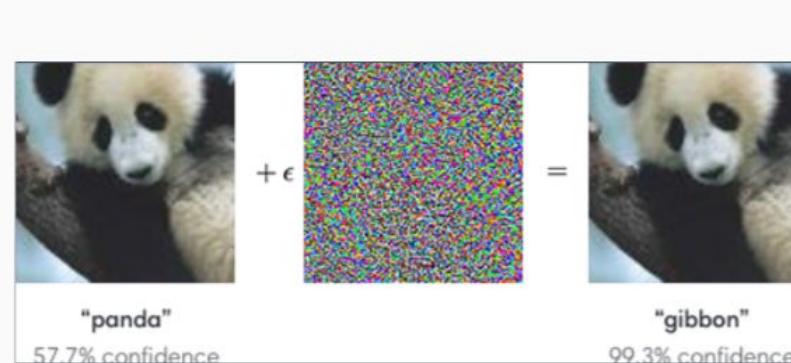
<https://www.quantamagazine.org/how-artificial-intelligence-is-changing-science-20190311/>

# Why get to the basics of ANNs?

## The Paradox of Deep Learning

Why does deep learning work so well despite

- high capacity (susceptible to overfitting)
- numerical instability (vanishing/exploding gradients)
- sharp minima (leading to overfitting)
- non-robustness (see figure)



No clear answers yet but ...

Slowly but steadily there is increasing emphasis on  
explainability and theoretical justifications!\*

Hopefully this will bring sanity to the proceedings !

---

\*<https://arxiv.org/pdf/1710.05468.pdf>

We would like to understand the key idea and algorithm that underlies almost all ANN/DL models.

# Outline for Module M9

- M9. Neural Networks
  - M9.0 Introduction/Motivation
  - **M9.1 Feed-forward neural networks**
    - **(Key Idea: Adaptive Basis functions and its Generalizations)**
    - (Network architecture, Network training & Backpropagation algo. sketch)
  - M9.2 Mixture Density Networks
  - M9.3 Concluding thoughts

# Key Idea: Adaptive Basis fns (aka Feature/Representation Learning)

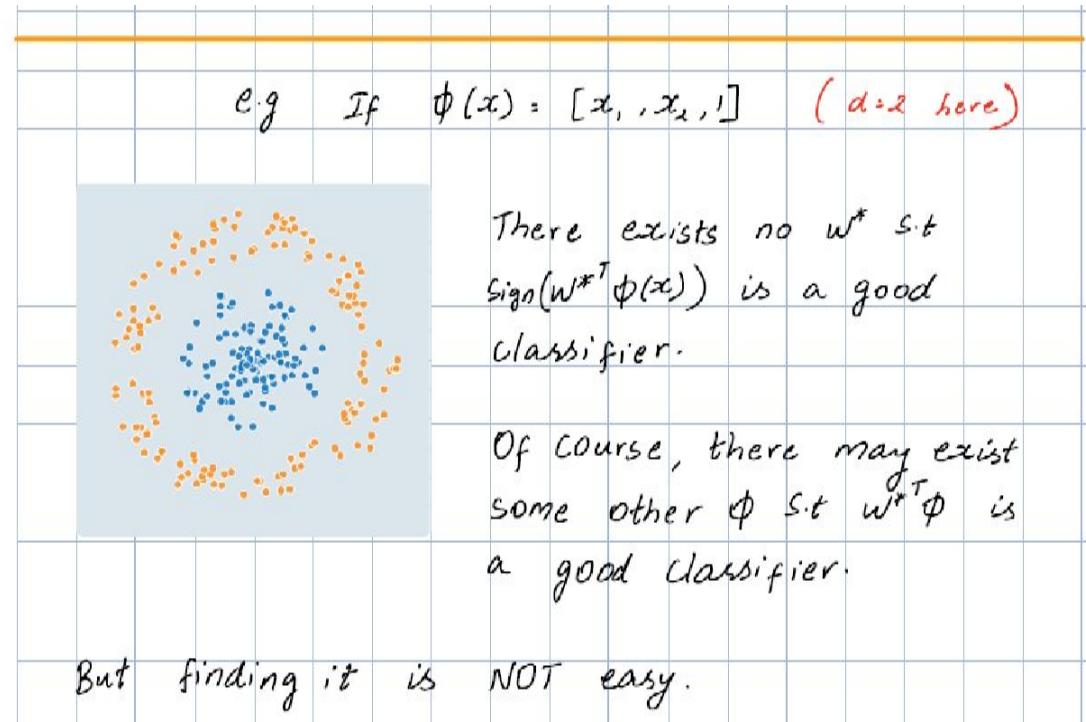
Sec 2

Beyond Linear models:

Linear Parameterisation :  $f(x) = w^T \phi(x)$  for some fixed  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^k$ .

All algs above make sense only if there exists a  $w^*$  s.t  
 $\text{sign}(w^{*T} \phi(x))$  is a good classifier.

This may not always be the case.



# ANNs: Key Idea

Key idea: What if we parameterise the feature mapping  $\phi$  also, and learn it along with the weight vector  $w$ ?

$$u \rightarrow d, \quad w \rightarrow d,$$

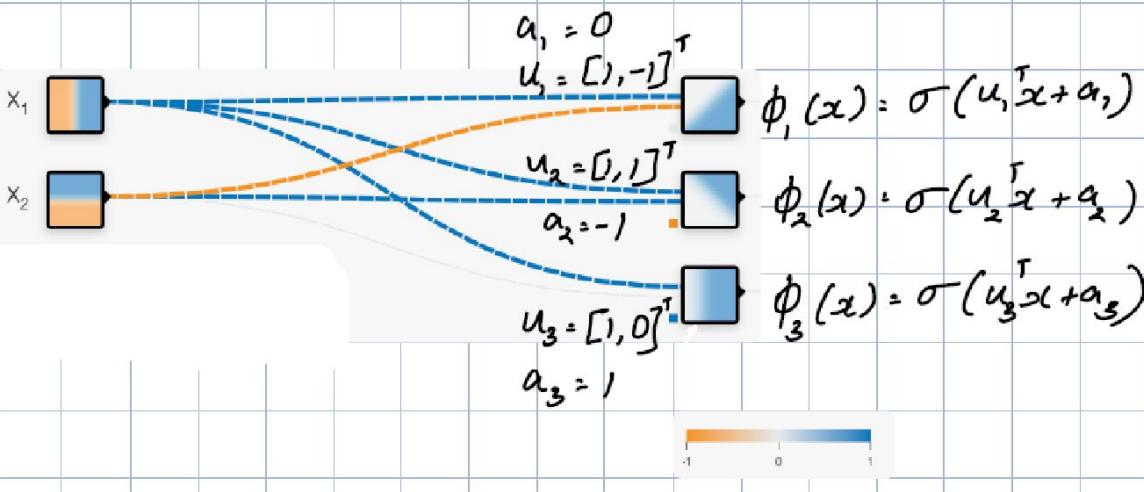
For example:  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^d$ ,

$$\phi(x) = \begin{bmatrix} \sigma(u_1^T x + a_1) \\ \sigma(u_2^T x + a_2) \\ \vdots \\ \sigma(u_d^T x + a_d) \end{bmatrix}$$

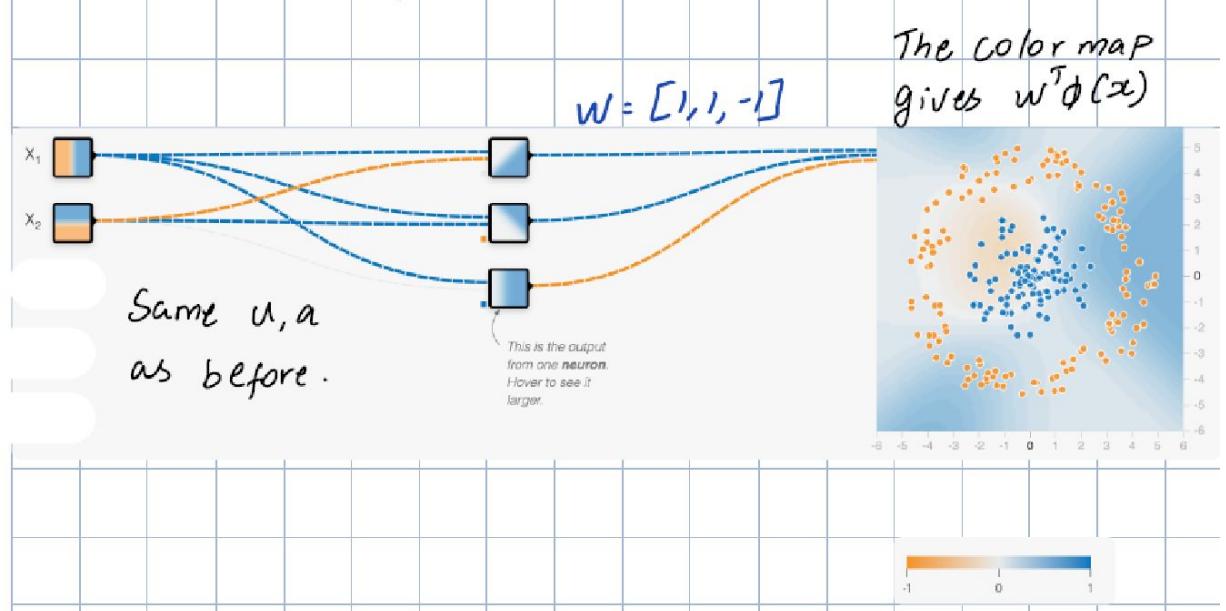
$\sigma: \mathbb{R} \rightarrow \mathbb{R}$  is a non-linear fn.  
Popular choices:  
 $\sigma(t) = \frac{1}{1+e^{-t}}$   
 $\sigma(t) = \tanh(t)$   
 $\sigma(t) = \max(0, t)$

# ANNs: Key Idea (example – specific u,a,w)

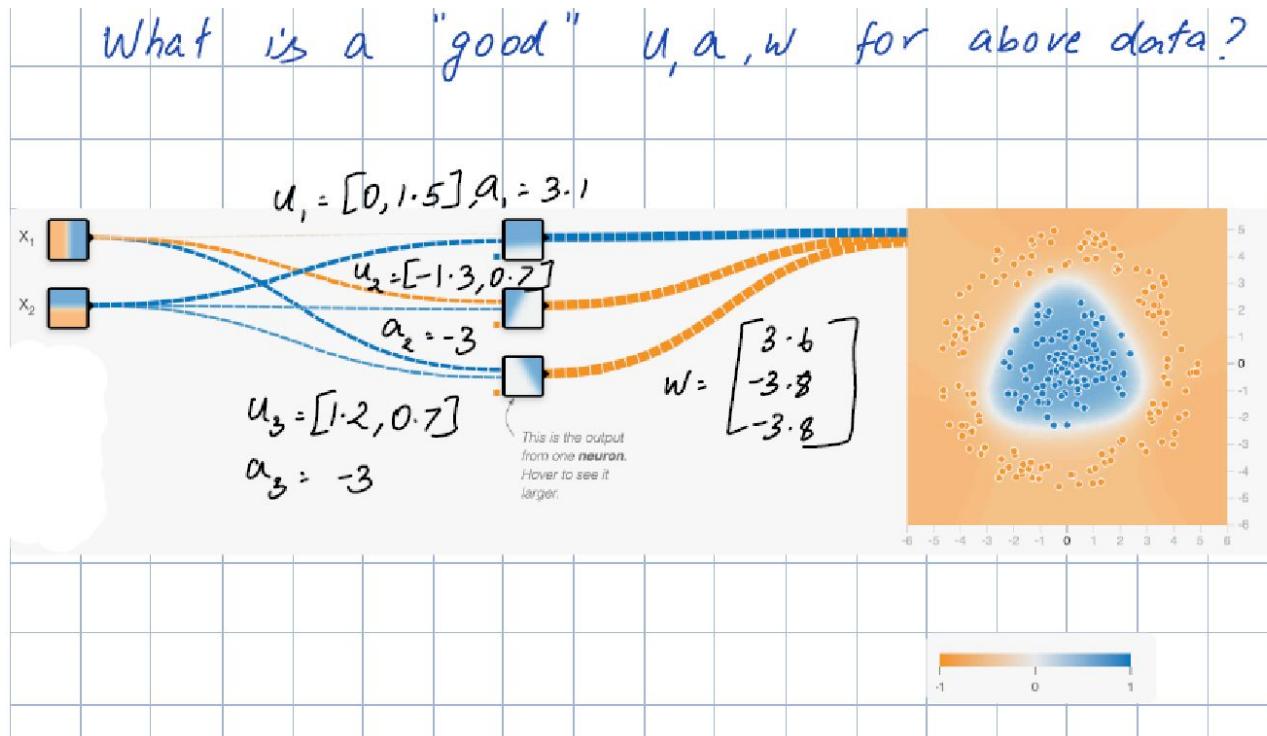
How does  $\phi$  look like for some  $u, a$ ?



How does  $w^T \phi(x)$  look like for some  $u, a, w$ ?



# ANNs: Key Idea (example – good u,a,w)



# Network training (getting a good $u, a, w$ ) helps us realize the key idea of feature/repn. learning!

How do we get a "good"  $u, a, w$ ?

- Simple gradient descent !!
- Efficient gradient computation: Backpropagation algorithm, an application of the chain rule for derivatives.

$$L(u, a, w) = \sum_{i=1}^n \log \left[ 1 + \exp(-y_i (w^T \phi(x_i))) \right]$$

$$\text{where } \phi(x) = \sigma(Ux + a) \quad U = \begin{bmatrix} u_1^T \\ \vdots \\ u_d^T \end{bmatrix} \quad a = \begin{bmatrix} a_1 \\ \vdots \\ a_d \end{bmatrix}$$

# Generalizations of the key idea

## Generalisations:

- (i) Can add even more "layers" e.g

$$\phi(x) = \sigma(V\tau(x) + b)$$

$$\tau(x) = \sigma(U\delta(x) + a)$$

$$\delta(x) = \sigma(Tx + c)$$



- (ii) Can use other non-linearities, e.g. tanh, ReLU

- (iii) Can use more sophisticated opt. algos than simple GD.

- (iv) Can add extra structure to the weights for special inputs. e.g. CNN for images, RNN for speech.

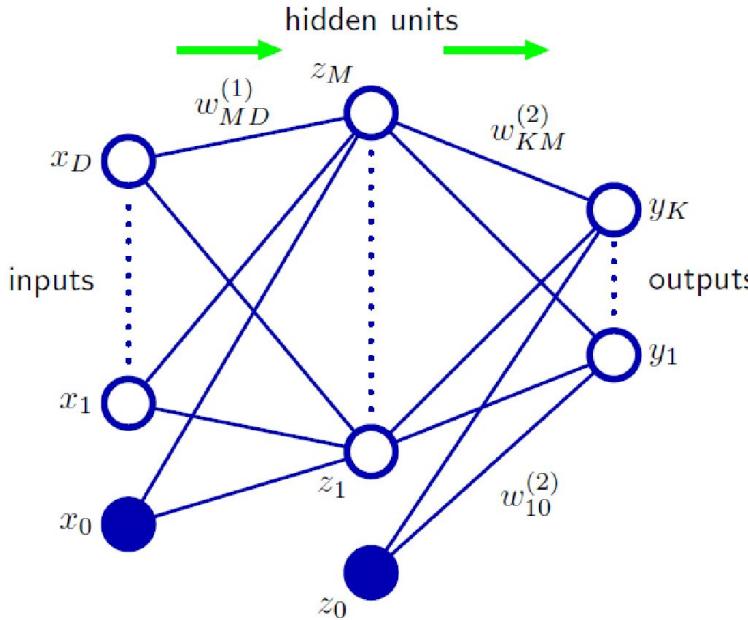
# Pros/cons of these ANN (generalizations)

Pros & Cons with neural nets:	
Pros:	Cons:
(i) More flexible than trying different fixed $\phi$ mappings.	(i) Optimisation need not always find a good param. even if it exists.
(ii) Allows to handle different types of input.	(ii) More tuning/validation of hyperparameters.
(iii) Works "well" in practice.	(iii) Lack of interpretability.

# Outline for Module M9

- M9. Neural Networks
  - M9.0 Introduction/Motivation
  - **M9.1 Feed-forward neural networks**
    - (Key Idea: Adaptive Basis functions and its Generalizations)
    - (**Network architecture, Network training & Backpropagation algo. sketch**)
  - M9.2 Mixture Density Networks
  - M9.3 Concluding thoughts

# Feed-forward Neural Network (FFN)



**Note change in notations:**

Training data:  $\{(x^{(n)}, t^{(n)})\}_{n=1 \text{ to } N'}$   
with  $x_n \in R^d$ ;

and  $t_n \in R^K$  for K regression problems,  
or  $t_n \in \{-1, +1\}^K$  for K two-class problems,  
or  $t_n \in \{\text{one. hot. encodings}\}$  for one  
multi(K)-class problem.

Predictions/outputs are:  $\{y_k\}_{k=1 \text{ to } K}$

Learning  $K$  models together permits sharing of  
weights of earlier layers across all  $K$  outputs!

Note FFNs aka Multilayer Perceptron -- but a misnomer as it's not multiple layers of perceptrons (which are non-continuous, non-differentiable step functions), but multiple layers of continuous non-linear functions like logistic functions!

[CMB]

# Output layer - popular choices

Problem	Activation fn. for output	
(1) Regression (K regression problems)		
(2) Classification (K two-class problems)		
(3) Classification (one K-class problem; K > 2)		

*MLE of a discriminative model is the same as minimizing the error/loss fn. above!*

[CMB]

# Exercise: Derivative for output layer

- Compute derivative  $\delta_k := \frac{\partial L_n}{\partial a_k} := \frac{\partial E_n}{\partial a_k}$  using chain rule for problems in previous page.

- Chain rule for:

- Problems (1),(2):

$$\frac{\partial E_n}{\partial a_k} = \frac{\partial E_n}{\partial y_k^{(n)}} \frac{\partial y_k^{(n)}}{\partial a_k} \quad \delta_k = y_k - t_k$$

- Problem (3):

$$\frac{\partial E_n}{\partial a_k} = \sum_{k'} \frac{\partial E_n}{\partial y_{k'}^{(n)}} \frac{\partial y_{k'}^{(n)}}{\partial a_k}$$

- Why are we interested in above?

$$\frac{\partial E_n}{\partial w_{ji}} = \frac{\partial E_n}{\partial a_j} \frac{\partial a_j}{\partial w_{ji}} = \delta_j z_i.$$

# Numerical differentiation wrt all weights

- If  $W$  is the total # of params., then below takes  $O(W)$  time for each  $w_{ji}$ , and hence  $O(W^2)$  for all params.

$$\frac{\partial E_n}{\partial w_{ji}} = \frac{E_n(w_{ji} + \epsilon) - E_n(w_{ji})}{\epsilon} + O(\epsilon)$$

$$\frac{\partial E_n}{\partial w_{ji}} = \frac{E_n(w_{ji} + \epsilon) - E_n(w_{ji} - \epsilon)}{2\epsilon} + O(\epsilon^2).$$

- Backprop. is simply chain rule realized via an  $O(W)$  Dynamic Programming (DP) algo.
  - essentially stores already computed derivatives wrt weights in later layers of the network to avoid redundant computations when computing derivatives wrt weights in earlier layers

# Backpropagation algorithm

## Error Backpropagation

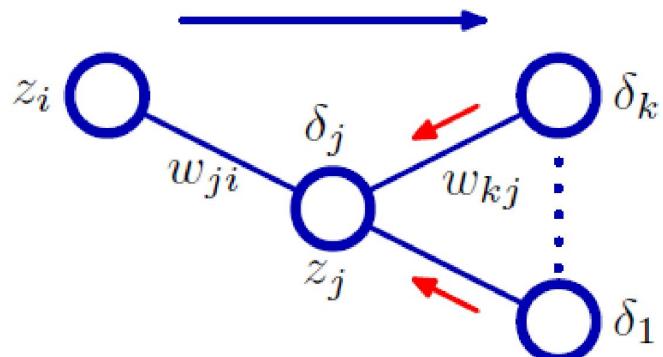
1. Apply an input vector  $\mathbf{x}_n$  to the network and forward propagate through the network using (5.48) and (5.49) to find the activations of all the hidden and output units.
2. Evaluate the  $\delta_k$  for all the output units using (5.54).
3. Backpropagate the  $\delta$ 's using (5.56) to obtain  $\delta_j$  for each hidden unit in the network.
4. Use (5.53) to evaluate the required derivatives.

$$a_j = \sum_i w_{ji} z_i \quad (5.48)$$

$$z_j = h(a_j). \quad (5.49)$$

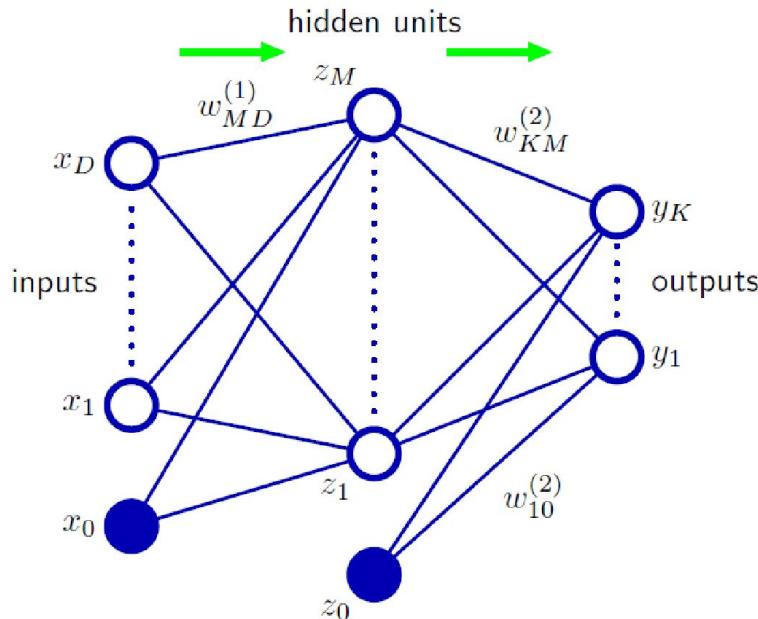
$$\delta_k = y_k - t_k \quad (5.54)$$

$$\frac{\partial E_n}{\partial w_{ji}} = \frac{\partial E_n}{\partial a_j} \frac{\partial a_j}{\partial w_{ji}} = \delta_j z_i. \quad (5.53)$$



$$\delta_j \equiv \frac{\partial E_n}{\partial a_j} = \sum_k \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial a_j} = h'(a_j) \sum_k w_{kj} \delta_k \quad (5.56)$$

# Exercise: When $\tanh(a)$ is the hidden unit activation fn, show that backprop. reduces to:



Next we compute the  $\delta$ 's for each output unit using

$$\delta_k = y_k - t_k. \quad (5.65)$$

Then we backpropagate these to obtain  $\delta$ s for the hidden units using

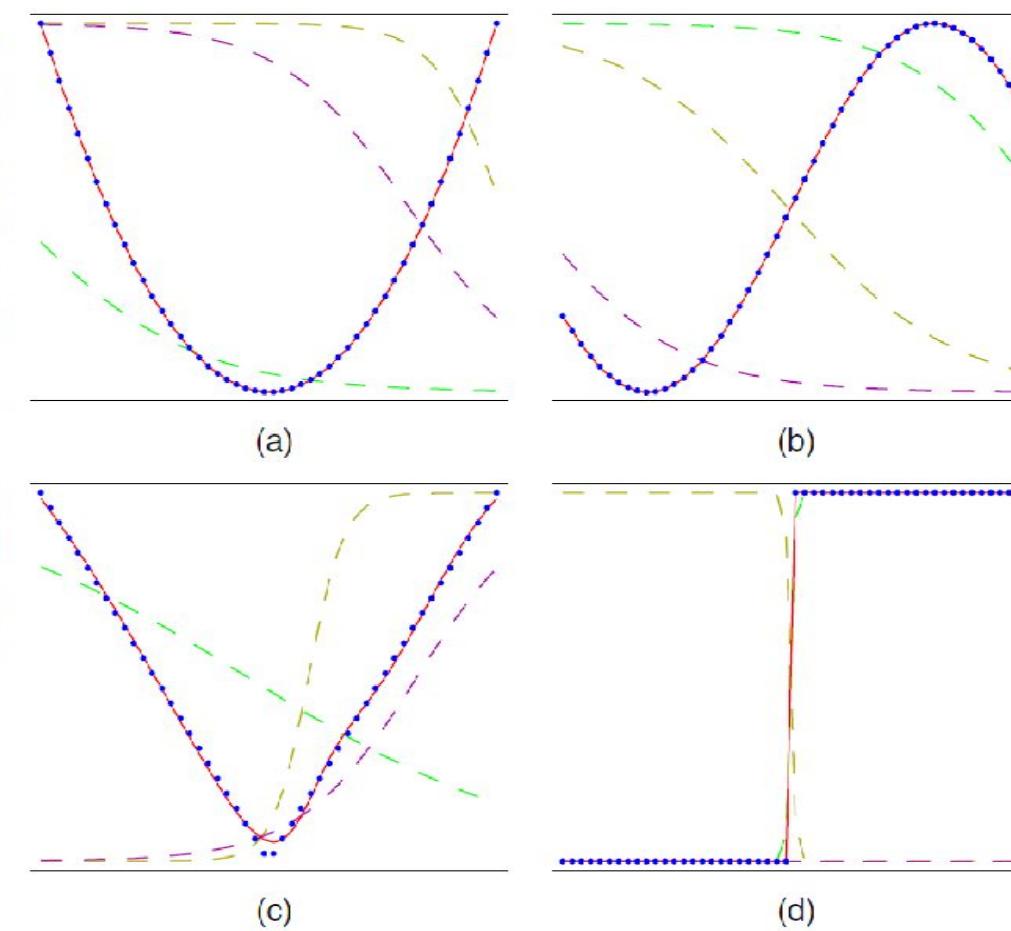
$$\delta_j = (1 - z_j^2) \sum_{k=1}^K w_{kj} \delta_k. \quad (5.66)$$

Finally, the derivatives with respect to the first-layer and second-layer weights are given by

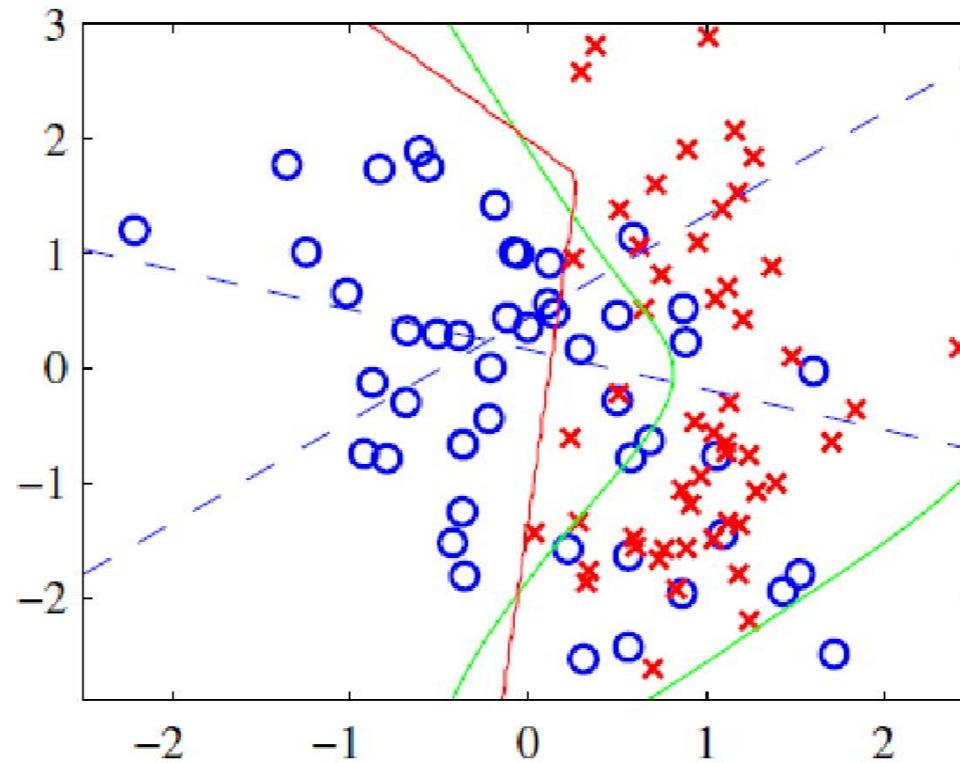
$$\frac{\partial E_n}{\partial w_{ji}^{(1)}} = \delta_j x_i, \quad \frac{\partial E_n}{\partial w_{kj}^{(2)}} = \delta_k z_j. \quad (5.67)$$

$$\tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}.$$

# What can you achieve with backprop + grad. descent? Example ANN for Regression



# Example ANN for Classification



# Numerical differentiation wrt all weights

- If  $W$  is the total # of params., then below takes  $O(W)$  time for each  $w_{ji}$ , and hence  $O(W^2)$  for all params.

$$\frac{\partial E_n}{\partial w_{ji}} = \frac{E_n(w_{ji} + \epsilon) - E_n(w_{ji})}{\epsilon} + O(\epsilon)$$

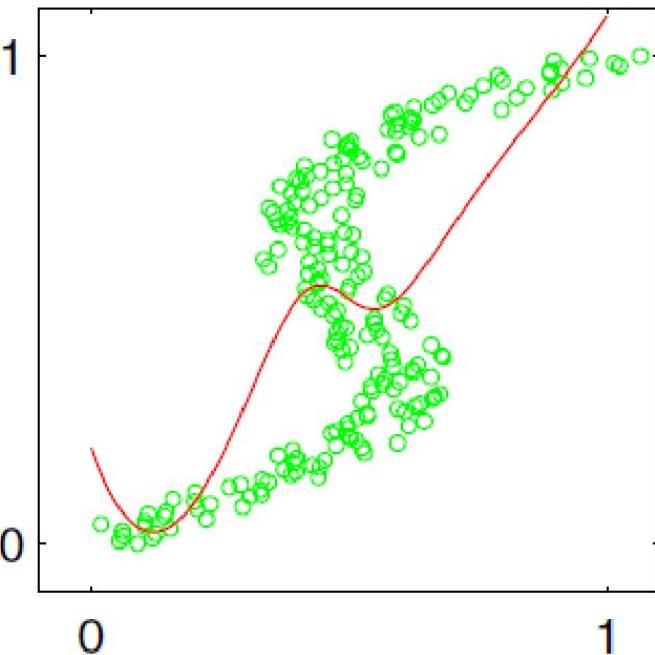
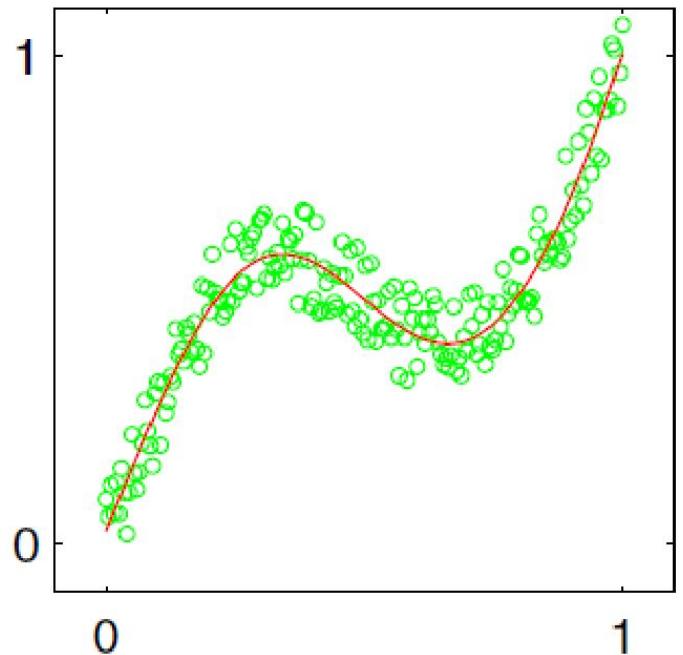
$$\frac{\partial E_n}{\partial w_{ji}} = \frac{E_n(w_{ji} + \epsilon) - E_n(w_{ji} - \epsilon)}{2\epsilon} + O(\epsilon^2).$$

- Backprop. is simply chain rule realized via an  $O(W)$  Dynamic Programming (DP) algo.
  - essentially stores already computed derivatives wrt weights in later layers of the network to avoid redundant computations when computing derivatives wrt weights in earlier layers

# Outline for Module M9

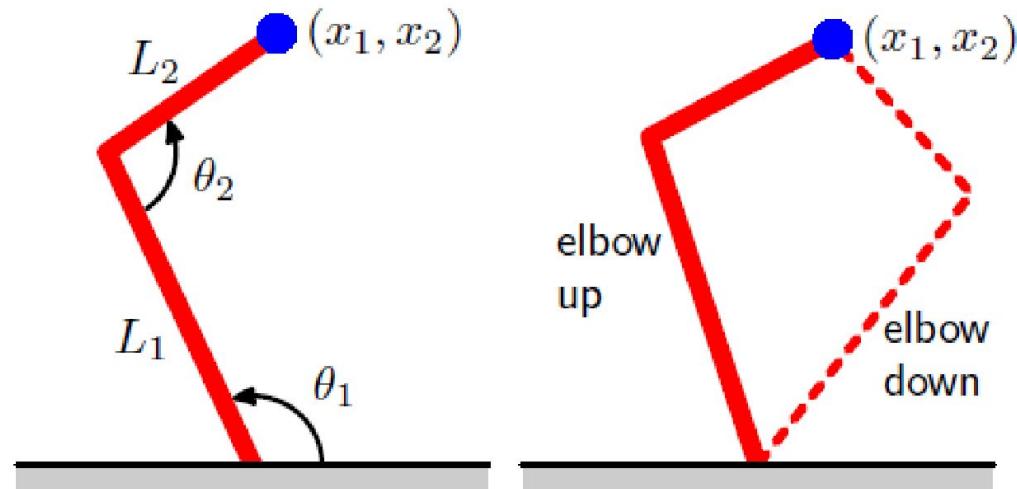
- M9. Neural Networks
  - M9.0 Introduction/Motivation
  - M9.1 Feed-forward neural networks
    - (Key Idea: Adaptive Basis functions and its Generalizations)
    - (Network architecture, Network training & Backpropagation algo. sketch)
  - **M9.2 Mixture Density Network**
  - M9.3 Concluding thoughts

Can ANN find a curve approximating  $t|x$  in both these datasets?



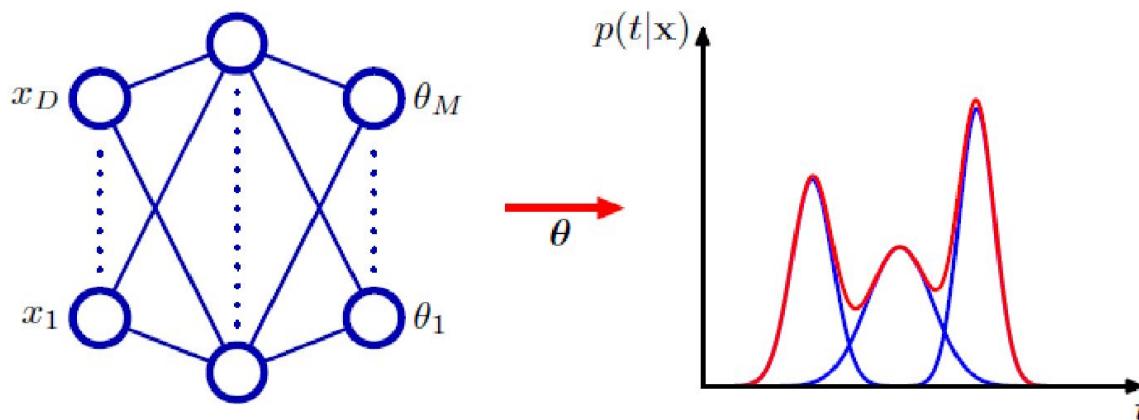
# Applications: robotics and disease prediction from symptoms

- Forward problem usually has a causal interpretation and a unique solution, whereas the inverse problem has multiple solutions/modes.



# Mixture density network (MDN): probabilistic model, network architecture & output unit choices

$$p(\mathbf{t}|\mathbf{x}) = \sum_{k=1}^K \pi_k(\mathbf{x}) \mathcal{N}(\mathbf{t}|\mu_k(\mathbf{x}), \sigma_k^2(\mathbf{x})).$$



$$\pi_k(\mathbf{x}) = \frac{\exp(a_k^\pi)}{\sum_{l=1}^K \exp(a_l^\pi)}.$$

$$\mu_k(x) = a_k^\mu \quad \sigma_k(\mathbf{x}) = \exp(a_k^\sigma).$$

[CMB]

# MDN Backprop algo. steps

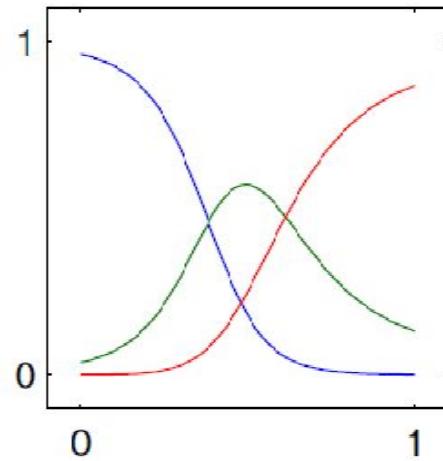
Step 1)  $E(\mathbf{w}) = - \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \mu_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \right\}$        $\gamma_k(\mathbf{t}|\mathbf{x}) = \frac{\pi_k \mathcal{N}_{nk}}{\sum_{l=1}^K \pi_l \mathcal{N}_{nl}}$

Step 2) Partial derivatives related to output layer:

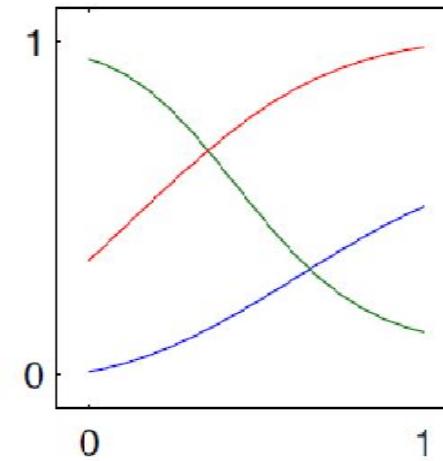
$$\frac{\partial E_n}{\partial a_k^\pi} = \pi_k - \gamma_k. \quad \frac{\partial E_n}{\partial a_{kl}^\mu} = \gamma_k \left\{ \frac{\mu_{kl} - t_l}{\sigma_k^2} \right\}. \quad \frac{\partial E_n}{\partial a_k^\sigma} = -\gamma_k \left\{ \frac{\|\mathbf{t} - \mu_k\|^2}{\sigma_k^3} - \frac{1}{\sigma_k} \right\}.$$

Steps 3,4) Propagate above derivatives backwards to compute gradient wrt all weights.

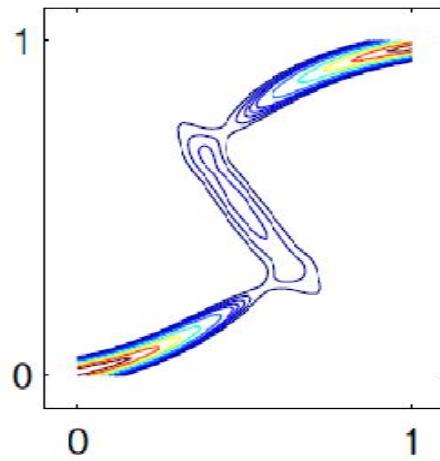
# Visualizing the MDN model



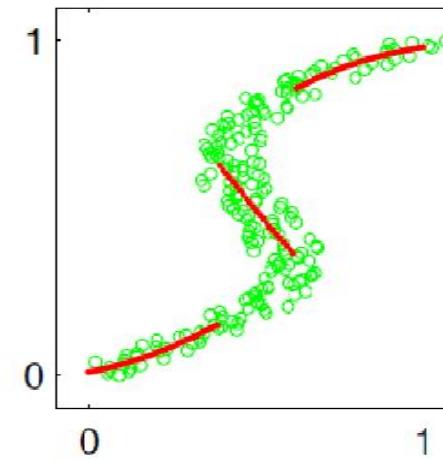
(a)



(b)



(c)



(d)

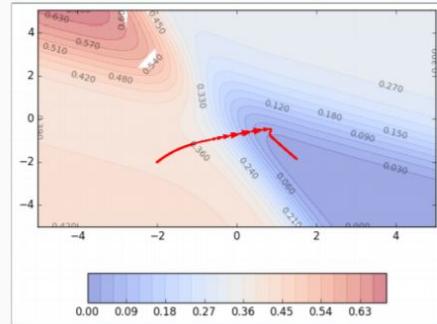
# Outline for Module M9

- M9. Neural Networks
  - M9.0 Introduction/Motivation
  - M9.1 Feed-forward neural networks
    - (Key Idea: Adaptive Basis functions and its Generalizations)
    - (Network architecture, Network training & Backpropagation algo. sketch)
  - M9.2 Mixture Density Network
  - **M9.3 Concluding thoughts**

# If you are interested to learn more...

## Better Optimization Methods

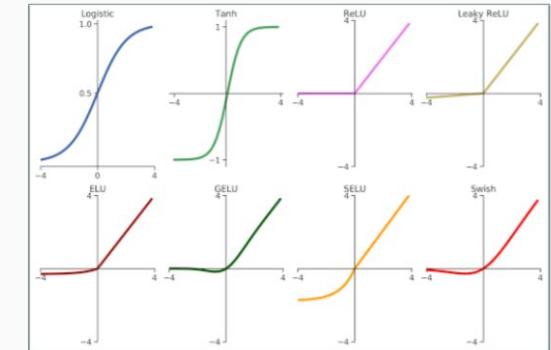
Faster convergence, better accuracies



## Better Activation Functions

We have come a long way from the initial days when the logistic function was the default activation function in NNs!

Over the past few years many new functions have been proposed leading to better convergence and/or performance!



Regularization of network (weights tying/sharing, etc.)

Calls for Sanity (Interpretable, Fair, Responsible, Green AI)

# Concluding thoughts & next steps

Non-linear method	(Non-linear) Basis functions	Objective function / OP
Vanilla extn. of linear models		Convex (unconstrained opt.)
SVM	<b>Selective</b> – center basis fns on training data points (dual/kernel view) and use training data to learn their weights and select a subset of them (non-zero weight support vectors) for eventual predictions.	Convex (constrained opt.)
Neural networks	<b>Adaptive</b> – Fix # of basis fns in advance, but allow them to be adaptive; parameterize basis fns and learn these parameters using training data.	Non-convex

**Next steps:** MDN is one way to probabilistically (soft-ly) combine different regression models (aka conditional mixtures), but there are also other complementary combined models / ensemble methods!

Thank you!