

# Lecture 24: (A Very Brief) Introduction to Artificial Neural Networks

MEGR 7080/8090: Dynamic System Learning and Estimation

Artur Wolek

Department of Mechanical Engineering and Engineering Science  
UNC Charlotte

Fall 2022

# Neural Networks (NNs) for Regression: Problem Setup

- Consider a static non-linear mapping from input  $\mathbf{x} \in \mathbb{R}^n$  to output  $\mathbf{y} \in \mathbb{R}^m$

$$\mathbf{y} = \mathbf{f}(\mathbf{x}) \quad (1)$$

- Regression NN Goal: learn (1) given *training data* (input/output pairs):

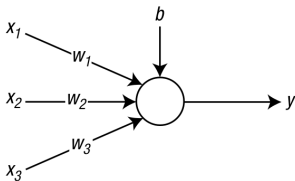
$$\mathcal{T} = \{(\mathbf{x}_1, \mathbf{t}_1), \dots, (\mathbf{x}_N, \mathbf{t}_N)\} \quad (2)$$

where  $\mathbf{x}_i$  is an input with corresponding (correct) output  $\mathbf{t}_i$ , and  $N$  is the number of training data points.

- Data may be synthetic, experimental, etc.

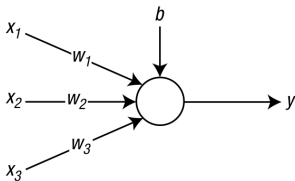
# Architecture: Node

- Node: receives inputs (vector) and produces a single output (scalar)
- Consider three inputs  $\mathbf{x} = [x_1, x_2, x_3]^T$  and scalar output  $y$ .
- The simplest one is a *linear transformation*: multiply inputs by weights  $w_1, w_2, w_3$ , respectively, and the result is summed and biased by a factor  $b$ :



# Architecture: Node

- Node: receives inputs (vector) and produces a single output (scalar)
- Consider three inputs  $\mathbf{x} = [x_1, x_2, x_3]^T$  and scalar output  $y$ .
- The simplest one is a *linear transformation*: multiply inputs by weights  $w_1, w_2, w_3$ , respectively, and the result is summed and biased by a factor  $b$ :



In matrix form:

$$y = \begin{bmatrix} w_1 & w_2 & w_3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + b = \mathbf{Ax} + b \quad (3)$$

# Nodes

- If instead there are two node, then there will be two outputs  $\mathbf{y} = [y_1, y_2]^T$  then the transformation that describes both nodes would be

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} w_1 & w_2 & w_3 \\ w_4 & w_5 & w_6 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \quad (4)$$

- If instead there are two node, then there will be two outputs  $\mathbf{y} = [y_1, y_2]^T$  then the transformation that describes both nodes would be

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} w_1 & w_2 & w_3 \\ w_4 & w_5 & w_6 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \quad (4)$$

- In general, any such transformation can be represented by the matrix equation

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b} \quad (5)$$

also called a *affine transformation*.

- The arrangement of inputs/nodes/outputs is arranged as set of *layers*.
- In the above example, the layer is parameterized by  $(\mathbf{A}, \mathbf{b})$

- Terminology: input layer, output layer, and intermediate (hidden) layers
- Each column of nodes is called a layer
- The output from one layer becomes the input to the next.

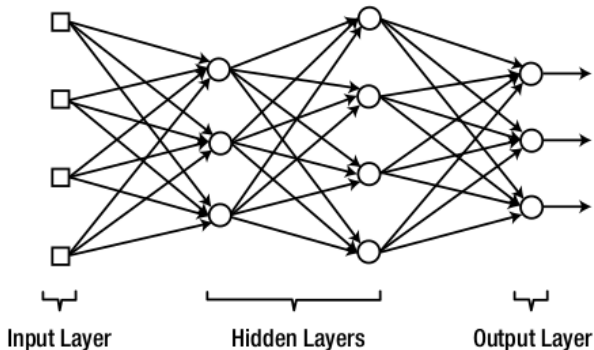
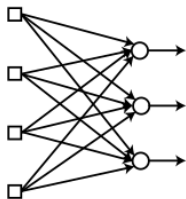
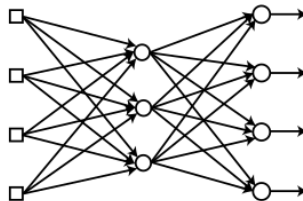


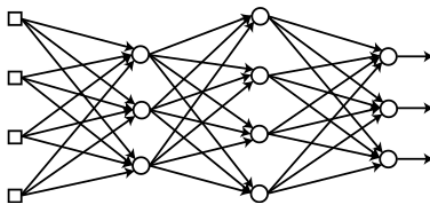
Image source: Kim 2017



Single-layer Neural Network



(Shallow) Multi-layer Neural Network



Deep Neural Network

Image source: Kim 2017



# Notation used here

- Consider a NN with  $d$  layers (including input/output layer)
- Input to each layer as  $\mathbf{q}_i$  where  $i \in \{1, \dots, d\}$ .

# Notation used here

- Consider a NN with  $d$  layers (including input/output layer)
- Input to each layer as  $\mathbf{q}_i$  where  $i \in \{1, \dots, d\}$ .
- Input is the first layer  $\mathbf{x} = \mathbf{q}_1 \in \mathbb{R}^n$
- Output is the last layer  $\mathbf{y} = \mathbf{q}_d \in \mathbb{R}^m$ .
- The intermediate  $\mathbf{q}_i$ 's are temp. variables of intermediate layers
- Note: layers can change size

# Notation used here

- Consider a NN with  $d$  layers (including input/output layer)
- Input to each layer as  $\mathbf{q}_i$  where  $i \in \{1, \dots, d\}$ .
- Input is the first layer  $\mathbf{x} = \mathbf{q}_1 \in \mathbb{R}^n$
- Output is the last layer  $\mathbf{y} = \mathbf{q}_d \in \mathbb{R}^m$ .
- The intermediate  $\mathbf{q}_i$ 's are temp. variables of intermediate layers
- Note: layers can change size
- Transition from one layer to next (for simple linear model):

$$\mathbf{q}_{i+1} = \mathbf{A}_i \mathbf{q}_i + \mathbf{b}_i \quad (6)$$

|              |   |
|--------------|---|
| Input        | $\mathbf{x} = \mathbf{q}_1$   |
| Second layer | $\mathbf{q}_2 = \mathbf{A}_1 \mathbf{q}_1 + \mathbf{b}_1$<br>$= \mathbf{A}_1 \mathbf{x} + \mathbf{b}_1$   |
| Third layer  | $\mathbf{q}_3 = \mathbf{A}_2 \mathbf{q}_2 + \mathbf{b}_2$<br>$= \mathbf{A}_2 (\mathbf{A}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2$<br>$= \mathbf{A}_2 \mathbf{A}_1 \mathbf{x} + (\mathbf{A}_2 \mathbf{b}_1 + \mathbf{b}_2)$  |
| Fourth layer | $\mathbf{q}_4 = \mathbf{A}_3 \mathbf{q}_3 + \mathbf{b}_3$<br>$= \mathbf{A}_3 (\mathbf{A}_2 \mathbf{A}_1 \mathbf{x} + (\mathbf{A}_2 \mathbf{b}_1 + \mathbf{b}_2)) + \mathbf{b}_3$<br>$= \mathbf{A}_3 \mathbf{A}_2 \mathbf{A}_1 \mathbf{x} + (\mathbf{A}_3 \mathbf{A}_2 \mathbf{b}_1 + \mathbf{A}_3 \mathbf{b}_2 + \mathbf{b}_3)$   |
|              | $\vdots$  |
| Output       | $\mathbf{y} = \mathbf{q}_d = \mathbf{A}_{d-1} \mathbf{q}_{d-1} + \mathbf{b}_{d-1}$<br>$= (\mathbf{A}_{d-1} \cdots \mathbf{A}_2 \mathbf{A}_1) \mathbf{x} + (\mathbf{A}_{d-1} \cdots \mathbf{A}_2) \mathbf{b}_1 + (\mathbf{A}_{d-1} \cdots \mathbf{A}_3) \mathbf{b}_2 + \cdots + \mathbf{b}_{d-1}$<br>$= \prod_{i=1}^{d-1} \mathbf{A}_i \mathbf{x} + \sum_{i=1}^{d-1} \prod_{j=i}^{d-2} \mathbf{A}_j \mathbf{b}_i + \mathbf{b}_d$ |

Conclusion: The multi-layer *linear* network can be replaced with an equivalent single-layer network with weights and bias vectors given by

$$\tilde{\mathbf{A}} = \prod_{i=1}^d \mathbf{A}_i \quad \text{and} \quad \tilde{\mathbf{b}} = \sum_{i=1}^d \prod_{j=i}^{d-1} \mathbf{A}_j \mathbf{b}_i + \mathbf{b}_d$$

so that the output becomes  $\mathbf{y} = \mathbf{q}_{d+1} = \tilde{\mathbf{A}}\mathbf{x} + \tilde{\mathbf{b}}$ .

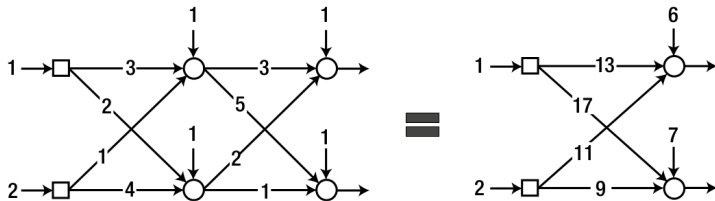


Image source: Kim 2017

# Activation Functions

- Output of each node passes through a nonlinear *activation function*  $\phi(\cdot)$ :

$$y_k = \phi(\mathbf{A}\mathbf{x} + b) \quad (7)$$

- Allows a richer set of models to be constructed.

$$y_k = \phi(\mathbf{A}\mathbf{x}) \quad (8)$$

# Activation Functions

- Output of each node passes through a nonlinear *activation function*  $\phi(\cdot)$ :

$$y_k = \phi(\mathbf{A}\mathbf{x} + b) \quad (7)$$

- Allows a richer set of models to be constructed.

$$y_k = \phi(\mathbf{A}\mathbf{x}) \quad (8)$$

- Note: Bias parameters can be absorbed into the set of weight parameters by defining an additional input variable  $x_0$  whose value is clamped at  $x_0 = 1$

# Activation Functions

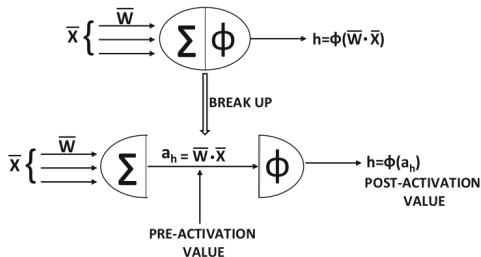
- Output of each node passes through a nonlinear *activation function*  $\phi(\cdot)$ :

$$y_k = \phi(\mathbf{A}\mathbf{x} + b) \quad (7)$$

- Allows a richer set of models to be constructed.

$$y_k = \phi(\mathbf{A}\mathbf{x}) \quad (8)$$

- Note: Bias parameters can be absorbed into the set of weight parameters by defining an additional input variable  $x_0$  whose value is clamped at  $x_0 = 1$





- Continuous or piece-wise differentiable (important for optimization)
- Examples:

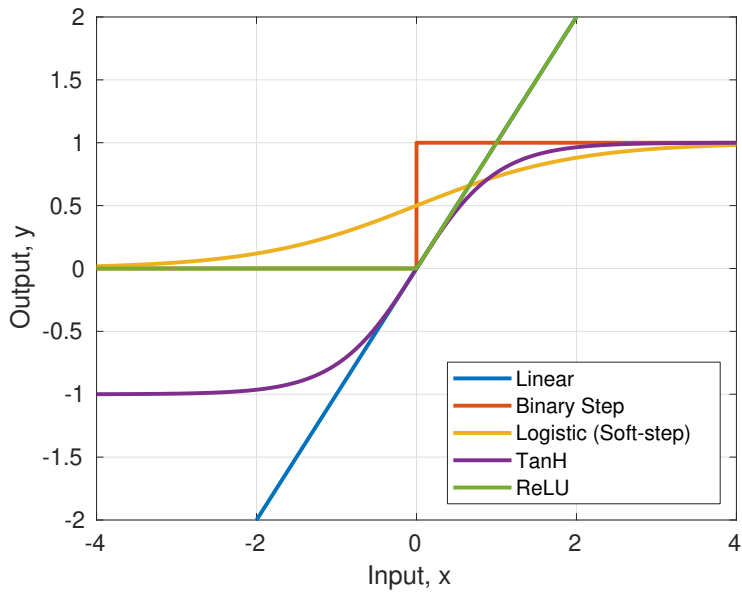
$$\phi(x) = x \quad \text{(linear)} \quad (9)$$

$$\phi(x) = \begin{cases} 0 & x \leq 0 \\ 1 & x > 0 \end{cases} \quad \text{(binary step)} \quad (10)$$

$$\phi(x) = \frac{1}{1 + \exp(-x)} \quad \text{(logistic (soft step))} \quad (11)$$

$$\phi(x) = \tanh(x) \quad \text{(TanH)} \quad (12)$$

$$\phi(x) = \begin{cases} 0 & x \leq 0 \\ x & x > 0 \end{cases} \quad \text{(rectified linear unit (ReLU))} \quad (13)$$



# Layer output with activation functions

The output vector at layer  $i + 1$  is determined by output at layer  $i$ :

$$\begin{bmatrix} (\mathbf{q}_{i+1})_1 \\ (\mathbf{q}_{i+1})_2 \\ \vdots \\ (\mathbf{q}_{i+1})_{N_{i+1}} \end{bmatrix} = \begin{bmatrix} \phi(\mathbf{A}_i^{(1)} \mathbf{q}_i) \\ \phi(\mathbf{A}_i^{(2)} \mathbf{q}_i) \\ \vdots \\ \phi(\mathbf{A}_i^{(N_i)} \mathbf{q}_i) \end{bmatrix} \quad (14)$$

Let  $\Phi : \mathbb{R}^{N_i} \rightarrow \mathbb{R}^{N_{i+1}}$ . We can represent the transition of the intermediate  $\mathbf{q}$  vectors through the affine transformations and activation functions as

$$\mathbf{q}_{i+1} = \Phi_i(\mathbf{q}_i) \quad (15)$$

The parameters of the  $i$ th layer are  $\mathbf{w}_i = \{w_{jk}\}$  for all  $j$  nodes in layer  $(i - 1)$  and  $k$  nodes in layer  $i$ . Weights determine  $\mathbf{A}_i^{(k)}$  matrices.

As before:

First (Input) Layer  $x = q_1$

Second Layer  $q_2 = \Phi_1(q_1)$

Third Layer  $q_3 = \Phi_2(q_2) = \Phi_3(\Phi_2(q_1))$

Fourth Layer  $q_4 = \Phi_4(q_3) = \Phi_4(\Phi_3(\Phi_3(q_1)))$

$\vdots$

Output Layer  $q_{d+1} = \Phi_d(q_d)$

$$\implies y = q_{d+1} = \Phi_d(\Phi_{d-1}(\cdots(\Phi_2(\Phi_1(x)))))$$

# Neural Network Weights

- The NN with nonlinear activation functions is

$$\begin{aligned}\hat{\mathbf{y}} &= \hat{\mathbf{f}}(\mathbf{x}; \mathbf{w}) \\ &= \Phi_d(\Phi_{d-1}(\cdots(\Phi_2(\Phi_1(\mathbf{x}; \mathbf{w}_1); \mathbf{w}_2) \cdots); \mathbf{w}_{d-1}); \mathbf{w}_d)\end{aligned}$$

- Passing in training data (2) through (??) gives predicted outputs:

$$\hat{\mathcal{T}} = \{(\mathbf{x}_1, \hat{\mathbf{y}}_1), \dots, (\mathbf{x}_N, \hat{\mathbf{y}}_N)\} . \quad (16)$$

- Ideally, predictions  $\{\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_N\}$  closely match correct outputs  $\{\mathbf{t}_1, \dots, \mathbf{t}_N\}$
- A common cost (loss) function  $J$  is a sum of squared errors

$$J(\mathbf{w}) = \sum_{i=1}^N \|\mathbf{y}_i - \hat{\mathbf{y}}_i\|_2 = \sum_{i=1}^N \|\mathbf{y}_i - \hat{\mathbf{f}}(\hat{\mathbf{x}}_i, \mathbf{w})\|_2 \quad (17)$$

where  $\|\cdot\|_2$  denotes the  $l^2$ -norm (i.e.,  $\|\mathbf{x}\|_2 = \sqrt{\sum_{k=1}^n |x_k|^2}$  for  $\mathbf{x} = [x_1, x_2, \dots, x_k]^T$ ).

# Training a Neural Network

- Training a NN involves adjusting the weights  $w$  to minimize loss (cost)
- Supervised learning: we know the correct outputs
- Once we have optimized the NN to fit the training data we can see how well it does on a *validation* data set that was not part of  $\mathcal{T}$
- Validation data can also be included in training

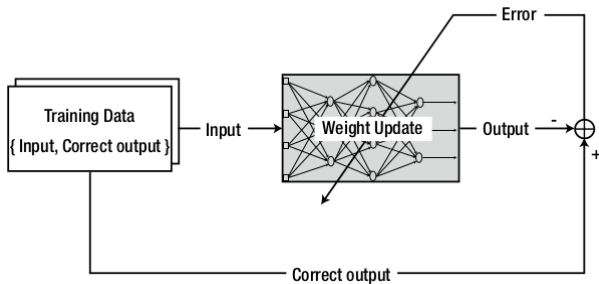


Image source: Kim 2017

# Gradient Descent

At each new iteration  $\tau + 1$ : update the parameters

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E(\mathbf{w}^{(\tau)}) \quad (18)$$

where  $\eta > 0$  is known as the learning rate.

**Figure 5.5** Geometrical view of the error function  $E(\mathbf{w})$  as a surface sitting over weight space. Point  $\mathbf{w}_A$  is a local minimum and  $\mathbf{w}_B$  is the global minimum. At any point  $\mathbf{w}_C$ , the local gradient of the error surface is given by the vector  $\nabla E$ .

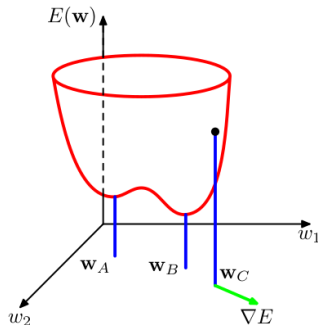


Image source: Bishop 2006

# Stochastic Gradient Descent

The “batch” error over all training data is

$$E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w}) \quad (19)$$

To speed up computation we can consider only a subset of the data points at each iteration.

- Sequential Gradient Descent: cycle through all data points
- Stochastic Gradient Descent: select data points at random

SGD can also help escape local minima



# How can we calculate the gradient of the error-function?

Consider error for just a single training data point

$$E_n = \frac{1}{2} \sum_{k=1}^{N_d} (y_k - t_k)^2 \quad (20)$$

where

- $n$ : identifier for the particular training data input
- $N_d$ : number of nodes in the output layer  $d$
- $y_k$ : predicted output for data input  $n$  at output node  $k$
- $t_k$ : correct output for data input  $n$  at output node  $k$

Note: predictions/outputs correspond to data index  $n$

The output  $q_k$  of an arbitrary node  $k$  in the network is a function of previous node outputs and weights

$$q_k = \phi \left( \sum_i^{\text{kth node's inputs}} w_{ki} q_i \right) \quad (21)$$

Let  $a_k$  denote the sum value prior to activation

$$a_k = \left( \sum_i^{\text{kth node's inputs}} w_{ki} q_i \right) \quad (22)$$

then

$$q_k = \phi(a_k) \quad (23)$$

Useful fact: Let

$$a_j = \left( \overset{j\text{th node's inputs}}{\sum_i} w_{ji} q_i \right) \quad \text{and} \quad a_k = \left( \overset{k\text{th node's inputs}}{\sum_m} w_{km} q_m \right) \quad (24)$$

Suppose output of node  $j$  feeds into node  $k$  for a particular value  $m = j$  so that

$$q_{m=j} = \phi(a_j) \quad (25)$$

No other  $q_m$  depend on  $a_j$  except for the case  $m = j$ . Then

$$\implies \frac{\partial a_k}{\partial a_j} = \frac{\partial a_k}{\partial q_j} \frac{\partial q_j}{\partial a_j} = w_{kj} \phi'(a_j) \quad (26)$$

The gradient of  $E_n$  with respect to a weight  $w_{ki}$

$$\frac{\partial E_n}{\partial w_{ki}} = \frac{1}{2} \frac{\partial}{\partial w_{ki}} \left( \underbrace{y_k}_{=\phi(a_k)} - t_k \right)^2 \quad (27)$$

$$= (y_k - t_k) \frac{\partial \phi}{\partial a_k} \frac{\partial a_k}{\partial w_{ki}} \quad (28)$$

$$= (y_k - t_k) \frac{\partial \phi}{\partial a_k} q_i \quad (29)$$
$$\underbrace{\hspace{10em}}_{\delta_k := \partial E_n / \partial a_k}$$

$$= \delta_k q_i \quad (30)$$

Delta terms are called “errors”

# Backpropagation

All of the  $q_i$  values are known because they were required to compute the output in the first “forward propagation” of the input

$$\frac{\partial E_n}{\partial w_{ki}} = \delta_k \underbrace{q_i}_{\text{from forward prop}} \quad (31)$$

The delta terms are easy to compute at the output layer

$$\delta_k := \frac{\partial E_n}{\partial a_k} = \left( \underbrace{y_k}_{\text{NN prediction}} - \underbrace{t_k}_{\text{correct output}} \right) \underbrace{\frac{\partial \phi}{\partial a_k}}_{\text{known derivative evaluated at } a_k} \quad (32)$$

The challenge is to find the delta terms at intermediate (hidden) nodes. Solution is known as *backpropagation*: use deltas at output to find deltas at previous layer and repeat.

Backpropagation argument:

- Variations in  $a_j$  give rise to variations in the error function  $E_n$  through variations in  $a_k$ s in the next layer (recall node  $j$  connects to several nodes  $k$  in the forward direction)

$$\delta_j := \frac{\partial E_n}{\partial a_j} = \sum_k^{\text{next layer}} \underbrace{\frac{\partial E_n}{\partial a_k}}_{\delta_k} \underbrace{\frac{\partial a_k}{\partial a_j}}_{w_{kj} \phi'(a_j)} \quad (33)$$

$$\implies \delta_j = \phi'(a_j) \sum_k^{\text{next layer}} w_{kj} \delta_k \quad (34)$$

Allows us to compute  $\delta_j$  from downstream deltas in the next layer

# Backpropagation Summary

- 1 Apply an input vector  $\mathbf{x}_n$  to the network and forward propagate to produce the outputs. Record the outputs  $q_k$  and summation values  $a_k$  of all hidden nodes.
- 2 Evaluate  $\delta_k$  for all of the output nodes (in the final layer)
- 3 Use backprop to find  $\delta_j$  for the second-to-last layer. Continue recursively working backwards one layer at a time to the input layer.
- 4 Assemble the error derivative  $\nabla E(\mathbf{w}^{(\tau)})$  from  $\frac{\partial E_n}{\partial w_{ki}} = \delta_k q_i$
- 5 Update weights
$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E(\mathbf{w}^{(\tau)}) \quad (35)$$
- 6 Go back to Step 1. Repeat.

When to stop?

# Overfitting

## Principle of Parsimony

“with competing theories or explanations, the simpler one, for example a model with fewer parameters, is to be preferred”

(Source: Wikipedia)

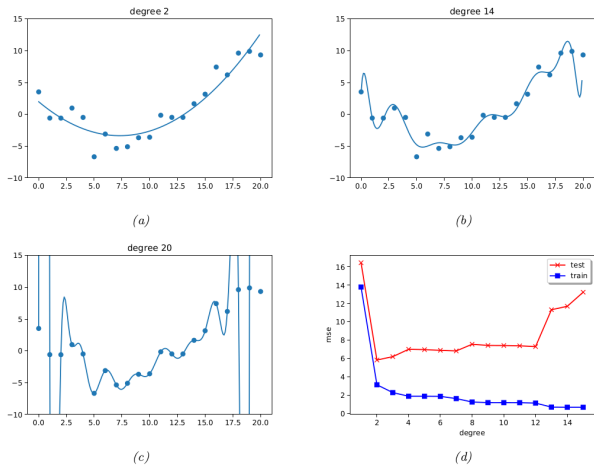


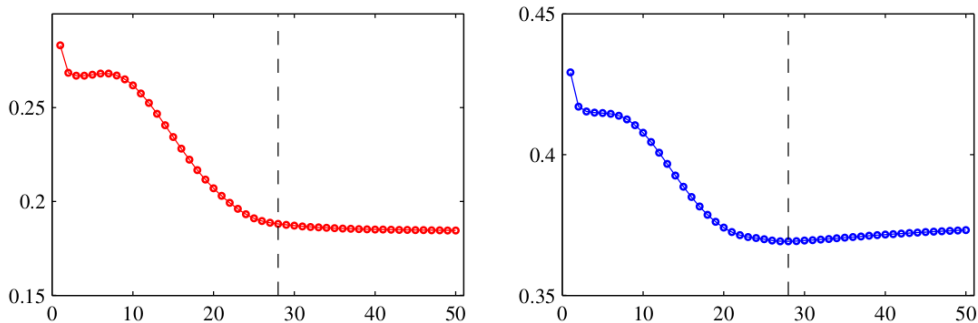
Figure 1.7: (a-c) Polynomials of degrees 2, 14 and 20 fit to 21 datapoints (the same data as in Figure 1.5). (d) MSE vs degree. Generated by [code.probl.ai/book1/1.7](https://code.probl.ai/book1/1.7).

Image source: Murphy 2022



# Simple approach: Early Stopping

- Stop training when validation error begins to increase



**Figure 5.12** An illustration of the behaviour of training set error (left) and validation set error (right) during a typical training session, as a function of the iteration step, for the sinusoidal data set. The goal of achieving the best generalization performance suggests that training should be stopped at the point shown by the vertical dashed lines, corresponding to the minimum of the validation set error.

Image source: Bishop 2006

# Training Issues

- Regularization tackles overfitting: add term  $\lambda C(\mathbf{w})$  to penalize model complexity where  $\lambda$  is a regularization parameter
- How to choose  $\lambda$ ?
  - Grid search (to maximize performance on validation data)
  - $k$ -folds cross-validation  $\rightarrow$  optimize  $\lambda^* \rightarrow$  train with all data  $\rightarrow \mathbf{w}^*$
- Another approach: dropout (randomly remove nodes during training)
- Tackling poor performance
  - Adjust network architecture
  - Adjust training parameters
  - Collect more data
  - Is the problem well posed?

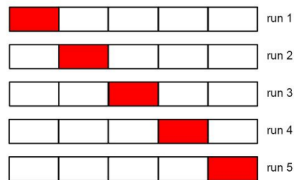


Figure 4.6: Schematic of 5-fold cross validation.

Image source: Murphy 2022

# Neural Network Implementation




|                           | Keras  | TensorFlow  | PyTorch  |
|---------------------------|---|--|---|
| Level of API              | high-level API <sup>1</sup>   | Both high & low level APIs   | Lower-level API <sup>2</sup>  |
| Speed                     | Slow  | High   | High  |
| Architecture              | Simple, more readable and concise   | Not very easy to use   | Complex <sup>3</sup>  |
| Debugging                 | No need to debug  | Difficult to debugging   | Good debugging capabilities   |
| Dataset Compatibility     | Slow & Small  | Fast speed & large   | Fast speed & large datasets   |
| Popularity Rank           | 1   | 2  | 3   |
| Uniqueness                | Multiple back-end support   | Object Detection Functionality   | Flexibility & Short Training Duration   |
| Created By                | Not a library on its own  | Created by Google  | Created by Facebook <sup>4</sup>  |
| Ease of use               | User-friendly   | Incomprehensive API  | Integrated with Python language   |
| Computational graphs used | Static graphs   | Static graphs  | Dynamic computation graphs <sup>5</sup>   |

Image Source: Yu-Cheng Kuo (Analytics Vidhya) [Link]

# Regression network in MATLAB

1) Randomly select initial conditions

```
x0Mat = ((rand([Nic 2])-0.5)*2)*3; % generate random ICs
```

2) Simulate the ODE, gather state-rate data

```
for i = 1:1:size(x0Mat,1)
    [t,X] = ode45(@(t,x) double_well(t,x,params), tspan, x0Mat(i,:) , options);
    tt = [tt; t];
    Xt = [Xt; X];
end
for i = 1:1:length(tt)
    dxdt = double_well(tt(i),Xt(i,:)',params);
    Xdot(i,1) = dxdt(1);
    Xdot(i,2) = dxdt(2);
end
```

3) Arrange into a table. Split into training/validation sets.

```
Xtrain = array2table([tt Xt Xdot], 'VariableNames',  
{ 't', 'x1', 'x2', 'x1dot', 'x2dot' });  
c = cvpartition(size(Xtrain,1), "Holdout", 0.30);  
trainingIndices = training(c);  
validationIndices = test(c);  
tblTrain = Xtrain(trainingIndices,:);  
tblValidation = Xtrain(validationIndices,:);
```

4) Train

```
x1dot_md1 = fitrnet(Xtrain, 'x1dot', 'PredictorNames', {'t', 'x1', 'x2'},  
"Standardize", true, "ValidationData", tblValidation, "LayerSizes", layerSize);
```

5) Predict (within ode45)

```
dxdt(1,1) = predict(x1dot_md1, [t xx'] );
```

# 10 ICs, 4 Layers, 30 nodes each, default settings

