



UMassAmherst
The Commonwealth's Flagship Campus

Neural Networks

Chaitra Gopalappa

Function Approximation/Fitting

- Widrow-Hoff works if functional form is known (model-based)
- What if function form is not know?
 - Use Neural Networks (model-free)

Neural networks

- Type of machine learning method
- Dates back to 1940's
- Nice read of the history 1940s to 2000:
<https://cs.stanford.edu/people/eroberts/courses/soco/projects/neural-networks/History/history1.html>
- Foundational work leading to current deep learning

Machine Learning Algorithms Cheat Sheet

```
graph TD
    START([START]) --> DR[Dimension Reduction]
    DR -- YES --> TM[Topic Modeling]
    DR -- NO --> HR[Have Responses]
    HR -- YES --> PN[Predicting Numeric]
    HR -- NO --> DP[Data Is Too Large]
    HR -- NO --> DR
    
    TM -- YES --> Prob[Probabilistic]
    Prob -- YES --> LDA[Latent Dirichlet Analysis]
    Prob -- NO --> SVD[Singular Value Decomposition]
    PC[Principal Component Analysis]
    
    PN -- YES --> SA[Speed or Accuracy]
    SA -- SPEED --> DT[Decision Tree]
    SA -- ACCURACY --> RF[Random Forest]
    RF --> NN[Neural Network]
    NN --> GB[Gradient Boosting Tree]
    
    DP -- YES --> NB1[Naïve Bayes]
    DP -- NO --> Exp[Explainable]
    Exp -- YES --> DT2[Decision Tree]
    Exp -- NO --> SA2[Speed or Accuracy]
    SA2 -- SPEED --> LS[Linear SVM]
    SA2 -- ACCURACY --> NB2[Naïve Bayes]
    SA2 -- ACCURACY --> LR[Logistic Regression]
    SA2 -- ACCURACY --> KSM[Kernel SVM]
    SA2 -- ACCURACY --> RF2[Random Forest]
    SA2 -- ACCURACY --> NN2[Neural Network]
    SA2 -- ACCURACY --> GB2[Gradient Boosting Tree]
    
    DP -- NO --> GP[Gaussian Mixture Model]
    GP --> PP[Prefer Probability]
    PP -- YES --> GMM[Gaussian Mixture Model]
    PP -- NO --> CV[Categorical Variables]
    CV -- YES --> KM[k-modes]
    CV -- NO --> PP
    CV -- YES --> HK[Hierarchical]
    HK -- YES --> HS[Hierarchical]
    HK -- NO --> NSK[Need to Specify k]
    NSK -- YES --> CV
    NSK -- NO --> DBSCAN[DBSCAN]
```

The flowchart is organized into four main sections, each with a colored header:

- Unsupervised Learning: Clustering** (Teal header): This section contains decision points for clustering algorithms. It starts with 'Prefer Probability' (leading to Gaussian Mixture Model if YES, or Categorical Variables if NO). 'Categorical Variables' leads to k-modes (if YES) or Hierarchical (if NO). 'Need to Specify k' leads to Hierarchical (if YES) or DBSCAN (if NO). 'Hierarchical' also leads to Hierarchical (if YES).
- Unsupervised Learning: Dimension Reduction** (Purple header): This section starts with 'Dimension Reduction'. If YES, it leads to 'Topic Modeling', which then leads to 'Probabilistic' (leading to Latent Dirichlet Analysis if YES, or Singular Value Decomposition if NO) or 'Principal Component Analysis' (if NO).
- Supervised Learning: Classification** (Yellow header): This section starts with 'Data Is Too Large'. If YES, it leads to Naïve Bayes. If NO, it leads to 'Explainable', which then leads to 'Speed or Accuracy'. 'Speed or Accuracy' leads to Linear SVM (if SPEED) or a choice between Naïve Bayes, Decision Tree, Logistic Regression, Kernel SVM, Random Forest, Neural Network, and Gradient Boosting Tree (if ACCURACY).
- Supervised Learning: Regression** (Pink header): This section starts with 'Predicting Numeric'. If YES, it leads to 'Speed or Accuracy', which then leads to Decision Tree (if SPEED) or a choice between Random Forest, Neural Network, and Gradient Boosting Tree (if ACCURACY). The 'Neural Network' option is highlighted with a white oval.

The flowchart begins with a 'START' node, which leads to 'Dimension Reduction'. If 'NO', it leads to 'Have Responses'. If 'NO' to 'Have Responses', it leads to 'Data Is Too Large'. If 'YES' to 'Have Responses', it leads to 'Predicting Numeric'. If 'NO' to 'Predicting Numeric', it leads to 'Data Is Too Large'. If 'YES' to 'Predicting Numeric', it leads to 'Speed or Accuracy'. If 'NO' to 'Speed or Accuracy', it leads to 'Data Is Too Large'. If 'YES' to 'Speed or Accuracy', it leads to the final algorithm selection based on 'SPEED' or 'ACCURACY'.

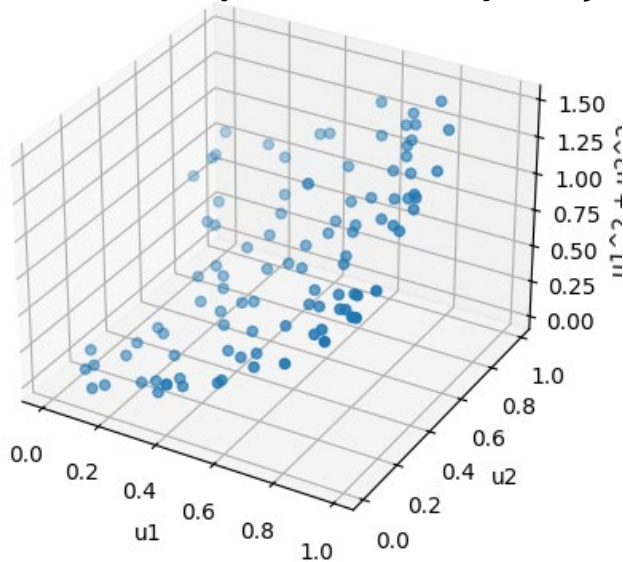
- In this class we will use Deep Learning/ Neural Network for Function Approximation (sometimes referred to as regression, though regression is one of the methods for function approximation)
- Deep Learning is a supervised Machine Learning method
 - Estimating a function that fits to the data

BASIC NN

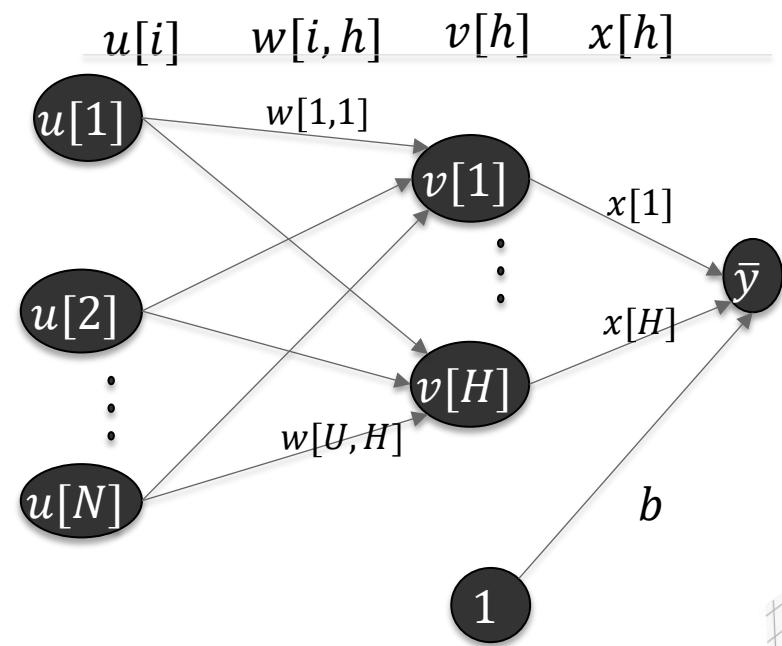
Nonlinear Neural Network (model- free):

- **Problem:** Fit a predictive model to the data. Suppose we do not know the analytical form of the function, but only know that $y = f(u_1, u_2, \dots)$

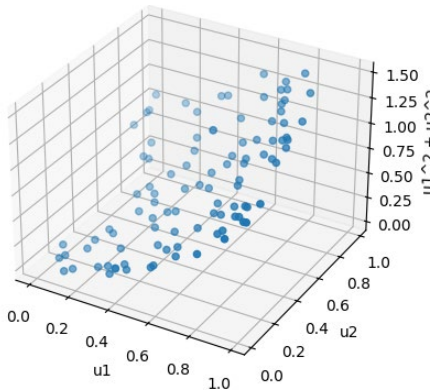
- Example $y = f(u_1, u_2)$



Neural Network (NN) architecture



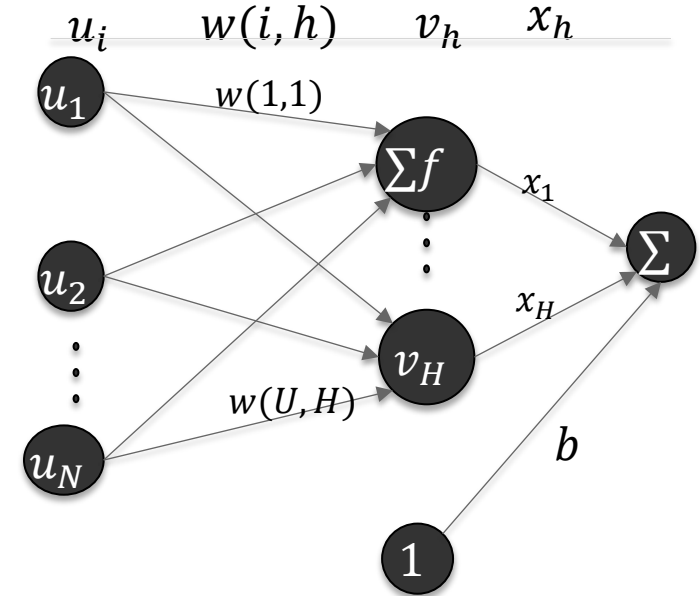
- $u[i]$ are the input nodes (independent variables of the function)
- $w[i, h]$ and $x[h]$ are the weights (equivalent to coefficients in polynomial equations) of the NN
 - $w[i, h]$: are the weights of the links from node $u[i]$ to node $v[h]$
 - $x[h]$: are the weights of the links from node $v[h]$ to \bar{y} (bar) to denote that it is an estimate of y)
- b is a bias node (equivalent to the intercept) term
- $v[h]$: are the nodes in the hidden layer and mathematically represented by **activation** functions (in below example we use a sigmoid function). **CORE of NN: adds non-linearity**



Neural Network terminologies

- Foundations are in linear algebra-All data are in arrays (vectors, matrices, 3D , ..)
- Most times notations not written in vector form (as on the right side figure)
- Will write in vector or matrix form on these slides to understand the foundations

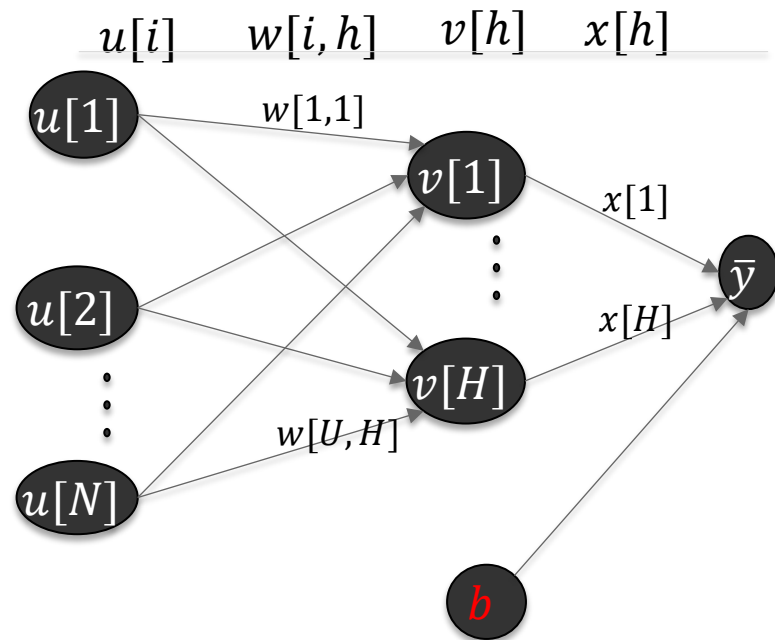
Suppose $y = f(\vec{u})$;
Corresponding NN architecture
for feedforward NN



Neural Network terminologies

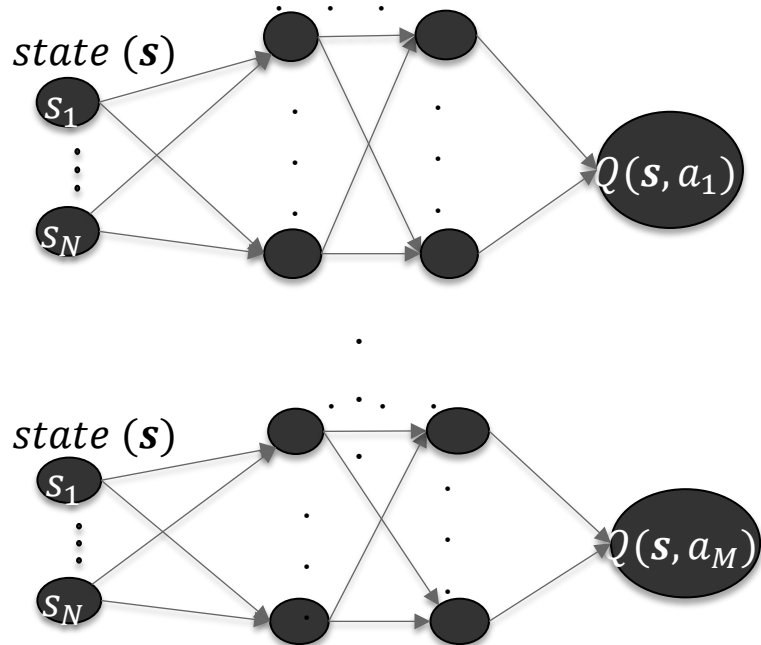
- NN architecture
 - Sometimes b written in the node, with nothing on arrow

Suppose $y = f(\vec{u})$;
Corresponding NN architecture
for feedforward NN

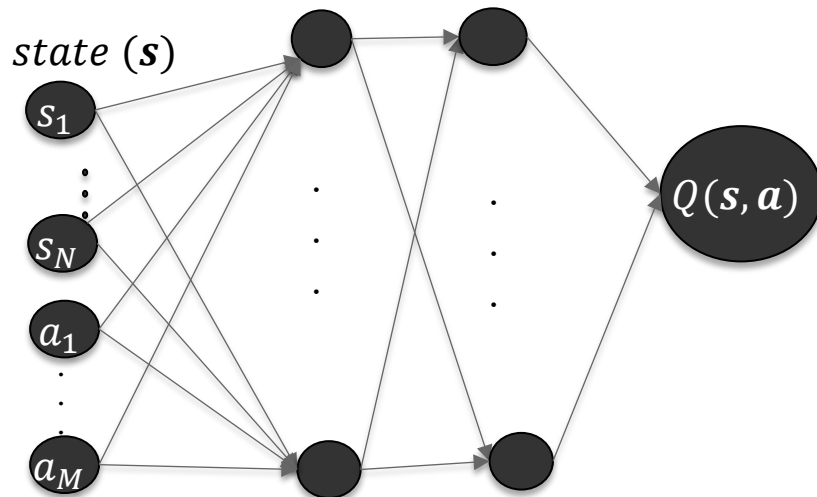


Example architectures - taking Q values

Example 1: One network for each action



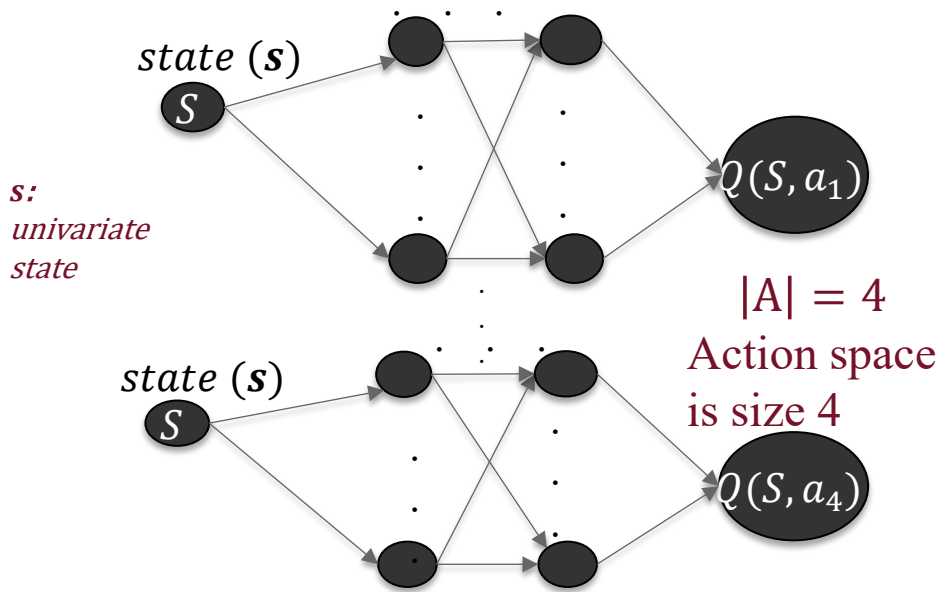
Example 2: Include action in input layer



s : state vector
 $s = [s_1, \dots, s_N]$

a : action vector
 $a = [a_1, \dots, a_M]$

Example of uni-variate state and action



Suppose,

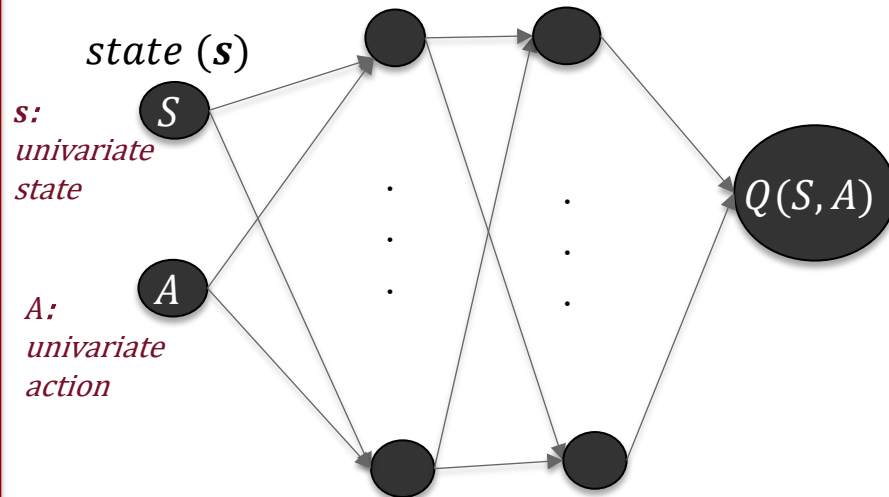
X_t = proportion of people with active infection

D_t = how often to test

State space: $S \in \mathbb{R}^1; S \in [0,1]$

Action space: $A = \{\text{test once a week, twice a week, three times a week, daily}\}$

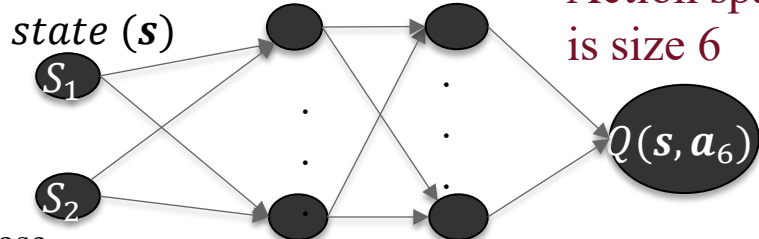
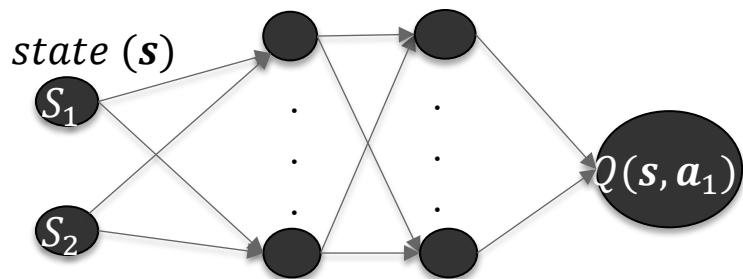
Include action in input layer



Action space: $A \in \mathbb{R}^1; A \in [1,30]$

Example of multi-variate state and action

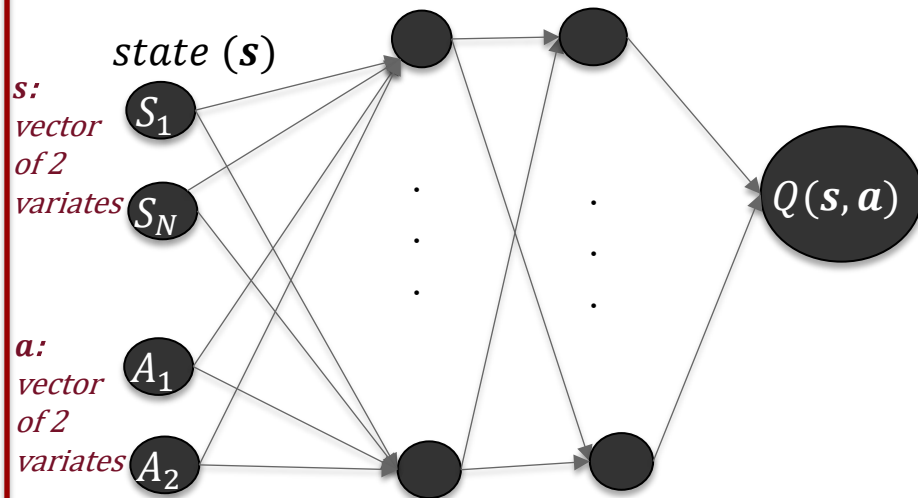
s:
vector
of 2
variables



$|A| = 6$
Action space
is size 6

Suppose,
 X_t = [proportion with active infection, proportion recovered]
 D_t = [how often to test in a week, what %lockdown]
 State space: $S \in \mathbb{R}^2$; (vector of 2 real random variables)
 Action space: $A = \{[\text{once}, 25\%], [\text{twice}, 25\%], [\text{thrice}, 25\%], [\text{once}, 50\%], [\text{twice}, 50\%], [\text{thrice}, 50\%]\}$

Include action in input layer



s:
vector
of 2
variables

a:
vector
of 2
variables

Action space: $a \in \mathbb{R}^2$; (vector of 2 real random variables)

s: state vector *a:* action vector
 $s = [S_1, \dots, S_N]$ $a = [A_1, \dots, A_M]$

Neural Network terminologies

- NN architecture

- Output layer (\bar{y} a vector) $\bar{y} = f(\vec{u})$
 - Figure shows univariate output (\bar{y})
- Input layer (\vec{u})
 - Input nodes ($u[i]$)
 - Number of nodes in input layer = $|\vec{u}|$
- Hidden layer/s (figure shows one hidden layer \vec{v})
- Weights (arrows; $w[.,.], x[.]$)
- Bias node (1 or b)
- Activation function/s (v)
 - Non-linear transformation applied to linear combination of $w[.,.], u[.]$

- Loss-function

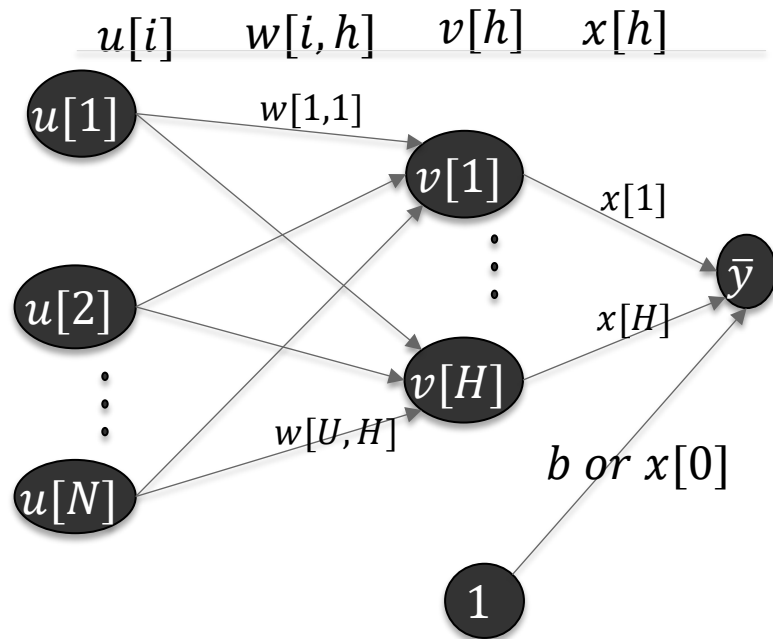
- Objective function (minimize loss/error)

- NN optimizer

- Algorithm to train the NN: solve for $w[.,.], x[.]$ to minimize loss

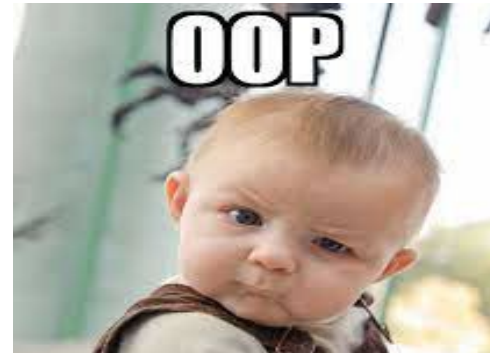
Suppose $y = f(\vec{u})$;

Corresponding NN architecture for feedforward NN



Guess what?

- **Pay attention to the variables;** they typically change (get into the habit of deviating from fixed variables)
- Here we are using u for independent variables (in W-H they were x);
- **Always define the variables for every problem,** do not expect the reader to guess; they always guess wrong



YOU ASKED ME TO GUESS;
AND I GUESSED WRONG

Neural Network for Function Fitting

- Objective for function fitting:

$$\min_{b, x, w} E = \min_{b, x, w} \sum_{p=1}^P (y_p - \bar{y}_p)^2$$

We can then apply Gradient Descent

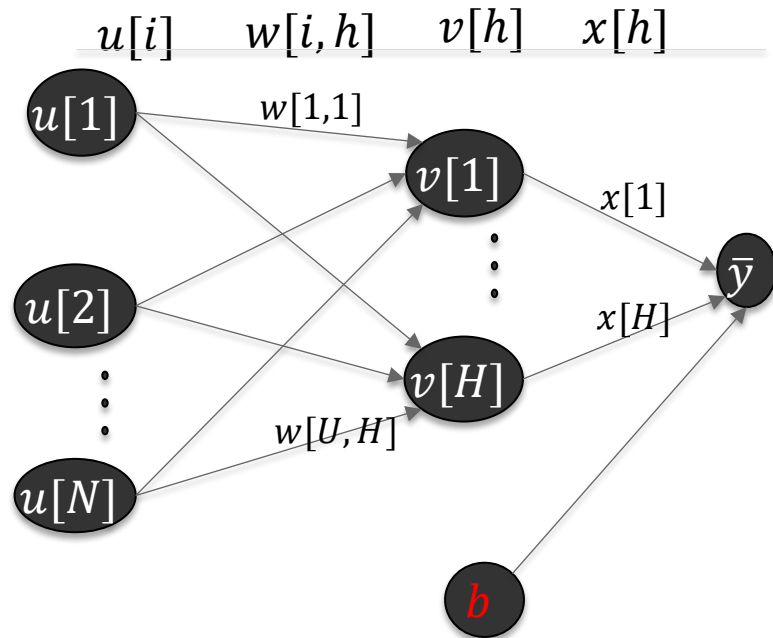
$$b \leftarrow b - \mu \frac{\partial(E)}{\partial b}$$

$$\omega[i, h] \leftarrow \omega[i, h] - \frac{\mu \partial(E)}{\partial \omega[i, h]}$$

$$x[h] \leftarrow x[h] - \mu \frac{\partial(E)}{\partial x[h]}$$

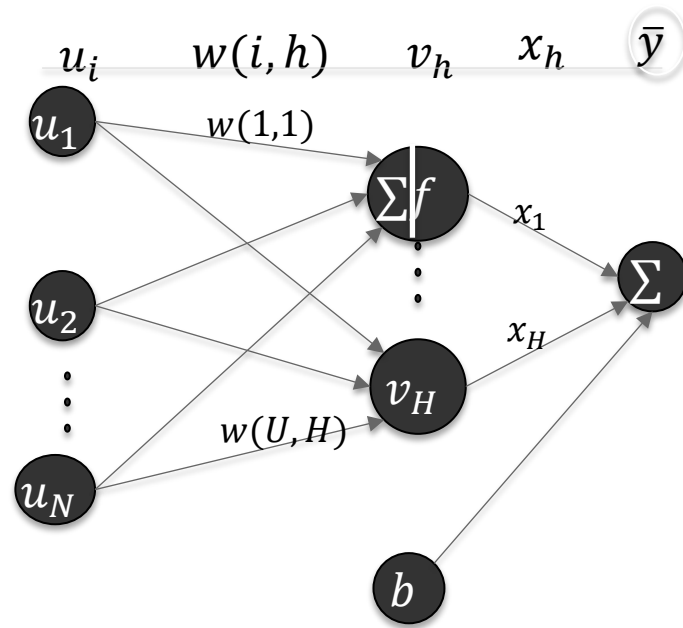
$$\bar{y} = ?$$

Suppose $y = f(\vec{u})$;



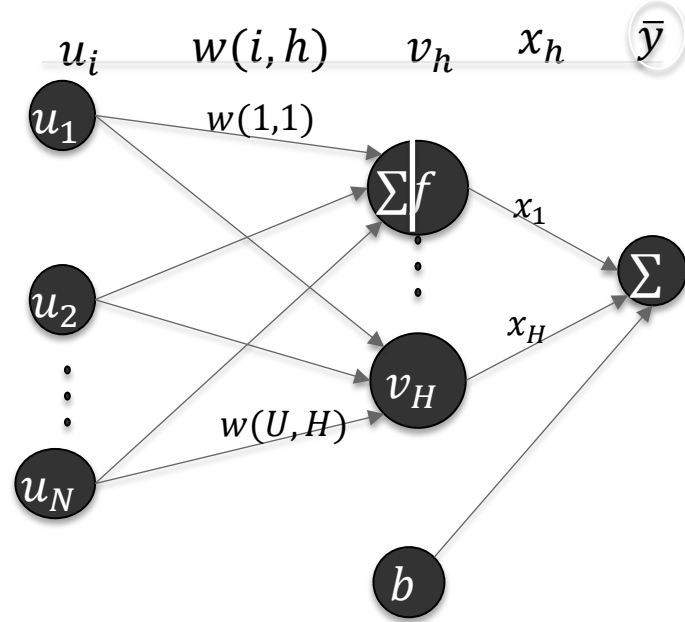
PROBLEM FORMULATION (mathematical representation of NN)

- We can mathematically represent the neural net with the following equations.
- $\bar{y} = b + \sum_{h=1,2,\dots,H} x[h]v[h]$
- $v[h]$ is a non-linear transformation applied to linear combination of $w[., h], u[.]$
- $v[h] = f(\sum_{i=1:N} w[i, h]u[i])$
- $v[h] = f(w[., h]u[.])$
- Suppose we use sigmoid activation function, then,
$$v[h] = \frac{1}{1+e^{-v^*[h]}}; \forall h, \text{ where}$$
$$v^*[h] = \sum_{i=1:N} w[i, h]u[i]$$



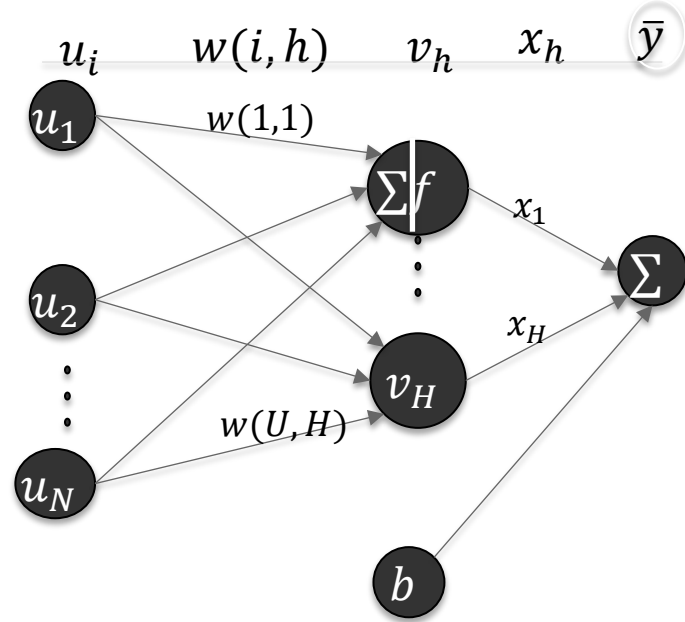
PROBLEM FORMULATION – with SIGMOID ACTIVATION

- $\bar{y} = b + \sum_{h=1,2,\dots,H} x[h]v[h]$
- $v[h] = \frac{1}{1+e^{-v^*[h]}}; \forall h$, where
- $v^*[h] = \sum_{i=1:N} w[i, h]u[i]$



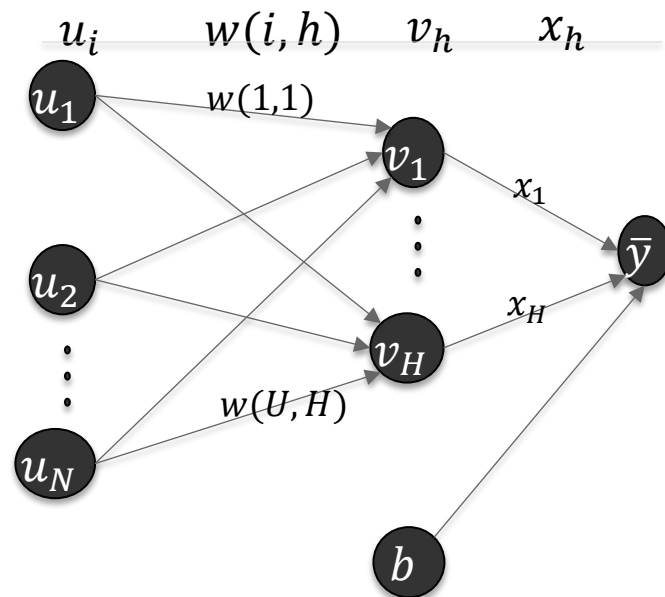
PROBLEM FORMULATION (SIGMOID ACTIVATION)

- $\bar{y} = b + \sum_{h=1,2,\dots,H} x[h]v[h]$
- $v[h] = \frac{1}{1+e^{-v^*[h]}}; \forall h$, where
- $v^*[h] = \sum_{i=1:N} w[i,h]u[i]$
- Objective is to $\min_{b,x,w} SSE = \sum_{p=1}^P (y_p - \bar{y}_p)^2$



SOLUTION ALGORITHM: Backprop (back propagation)

- **Backprop concepts**
- Objective is to $\min_{b,x,w} SSE \sim \min_{b,x,w} \frac{SSE}{2}$
 - $SSE = f(b, \omega[.,.], x[.])$ (Notice this doesn't have u because they are data samples)
 - $b, \omega[.,.]$, and $x[.]$ are the decision variables from perspective of Backprop.
 - \bar{y}_p = output from trained net for \vec{u}_p ; p are the data samples
 - y_p = actual data
 - P = number of samples; N = number of input nodes



Backprop algorithm

1. Initialize

1. Initialize $b, \omega(\cdot, \cdot)$, and $x(\cdot)$ to random values
2. Set SSE_{old} to very large value.
3. Set $m = 0$ (iteration number)

2. Compute $v_p^*[h] = \sum_i w[i, h] u_p[i]$ for each p and each h .

3. Compute $v_p[h] = \frac{1}{1 + e^{-v_p^*[h]}}$ for each p and each h

4. Compute $\bar{y}_p = b + \sum_h x[h] v_p[h]$ for each p (data samples)

5. Apply SD transformations ($E = \frac{1}{2} \sum_{p=1}^P (y_p - \bar{y}_p)^2$)

$$b_m \leftarrow b_{m-1} - \mu \frac{\partial(E)}{\partial b}$$

6. $\omega_m[i, h] \leftarrow \omega_{m-1}[i, h] - \mu \frac{\partial(E)}{\partial \omega[i, h]}$

$$x_m[h] \leftarrow x_{m-1}[h] - \mu \frac{\partial(E)}{\partial x[h]}$$

7. Set $m = m + 1$

1. Calculate $E_{new} = \sum_p (y_p - \bar{y}_p)^2$
2. Update μ
3. If $|E_{new} - E_{old}| < \text{tolerance} \rightarrow STOP$ Otherwise set $E_{old} = E_{new}$ got step 2.

Feed
forward

Recollect W-H for regression- uses Steepest Descent (SD)

• Initialize

- Set $w[i]$ to values between 0 and 1;
- Set E_{old} (the SSE) to a large number
- Set $m = 0$; Set μ

• Compute $\bar{y}_p = \sum_{i=0}^N \omega[i] x_p[i]$ for each p (data samples)

• Apply SD transformations ($E = \frac{1}{2} \sum_{p=1}^P (y_p - \bar{y}_p)^2$)

$$\omega_{m+1}[i] \leftarrow w_m[i] + \mu \frac{\partial(E)}{\partial w[i]}$$

$$\frac{\partial(E)}{\partial w[i]} = \sum_{p=1}^m (y_p - \bar{y}_p) x_p[i]$$

• Set $m = m + 1$.

- Calculate $E_{new} = \sum_p (y_p - \bar{y}_p)^2$
- Update μ .
- If $|E_{new} - E_{old}| < \text{tolerance}$ STOP. Otherwise set $E_{old} = E_{new}$ and go back to step 2.

Backprop algorithm

1. Initialize

1. Initialize $b, \omega(. , .)$, and $x(.)$ to random values
2. Set SSE_{old} to very large value.
3. Set $m = 0$ (iteration number)

2. Compute $v_p^*[h] = \sum_i w[i, h] u_p[i]$ for each p and each h .

3. Compute $v_p[h] = \frac{1}{1 + e^{-v_p^*[h]}}$ for each p and each h

4. Compute $\bar{y}_p = b + \sum_h x[h] v_p[h]$ for each p .

5. Apply SD transformations ($E = \frac{1}{2} \sum_{p=1}^P (y_p - \bar{y}_p)^2$)

$$b_m \leftarrow b_{m-1} - \mu \frac{\partial(E)}{\partial b}$$

6. $\omega_m[i, h] \leftarrow \omega_{m-1}[i, h] - \frac{\mu \partial(E)}{\partial \omega[i, h]}$

$$x_m[h] \leftarrow x_{m-1}[h] - \mu \frac{\partial(E)}{\partial x[h]}$$

7. Set $m = m + 1$

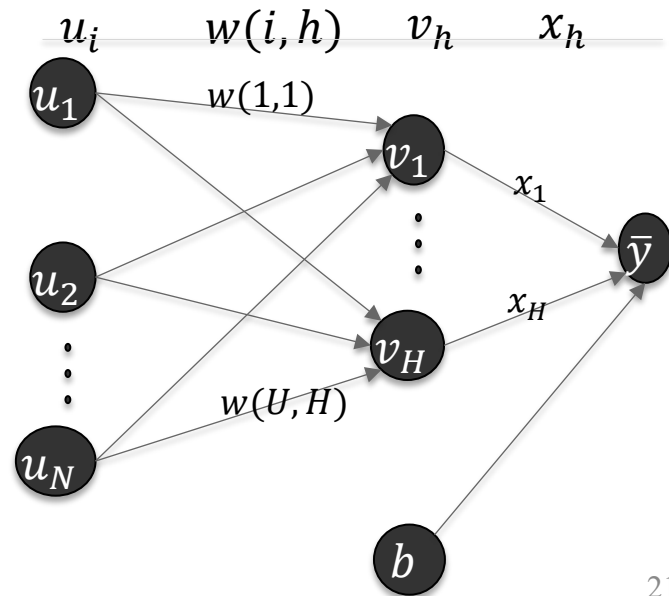
1. Calculate $E_{new} = \sum_p (y_p - \bar{y}_p)^2$
2. Update μ
3. If $[E_{new} - E_{old}] < \text{tolerance} \rightarrow STOP$ Otherwise set $E_{old} = E_{new}$ got step 2.

Feed
forward

- Do multiple iterations of the algorithm to start at different initial points
- Initialize $b, \omega[. , .]$, and $x[.]$ to select from a random range; Try a few different ranges
- Try a few different options for learning rate ($\mu = \frac{A}{B+m}; \mu = \frac{1}{m}; \mu = 0.001(\text{small constant})$)

Derivation of Backprop derivatives

- Main transformation is the steepest descent $\vec{x} \leftarrow \vec{x} - \mu \nabla f(\vec{x})$; but here we are solving for $\omega[.,.]$, and $x[.]$
- $b_m \leftarrow b_{m-1} - \mu \frac{\partial(E)}{\partial b}$
- $\omega_m[i, h] \leftarrow \omega_{m-1}[i, h] - \frac{\mu \partial(E)}{\partial \omega[i, h]}$
- $x_m[h] \leftarrow x_{m-1}[h] - \mu \frac{\partial(E)}{\partial x[h]}$



$$\frac{\partial E}{\partial x[h]}=? \quad \frac{\partial E}{\partial b}=?$$

$$\begin{aligned} \bullet \quad \frac{\partial E}{\partial x[h]} &= \frac{1}{2} \sum_{p=1}^P \frac{\partial}{\partial x[h]} (y_p - \bar{y}_p)^2 \\ &= \frac{2}{2} \sum_{p=1}^P (y_p - \bar{y}_p) \frac{\partial}{\partial x[h]} (y_p - \bar{y}_p) \\ &= \sum_p (y_p - \bar{y}_p) \left(-\frac{\partial \bar{y}_p}{\partial x[h]} \right) \\ &= \sum_p (y_p - \bar{y}_p) \left(-\frac{\partial (b + \sum_h x[h] v_p[h])}{\partial x[h]} \right) \end{aligned}$$

$$\frac{\partial E}{\partial x[h]} = -\sum_p (y_p - \bar{y}_p) v_p[h]$$

$$\bullet \quad \frac{\partial E}{\partial b} = -\sum_p (y_p - \bar{y}_p)$$

$$\frac{\partial E}{\partial \omega(i, h)} = ?$$

$$\begin{aligned} \frac{\partial E}{\partial \omega[i, h]} &= \frac{1}{2} \sum_p \frac{\partial}{\partial \omega[i, h]} (y_p - \bar{y}_p)^2 \\ &= \frac{1}{2} \sum_p 2(y_p - \bar{y}_p) \frac{\partial}{\partial \omega[i, h]} (y_p - \bar{y}_p) \\ &= \sum_p (y_p - \bar{y}_p) \left(-\frac{\partial}{\partial w[i, h]} \bar{y}_p \right) \\ &= - \sum_p (y_p - \bar{y}_p) \cdot \left(\frac{\partial \bar{y}_p}{\partial v_p[h]} \frac{\partial v_p[h]}{\partial v_p^*[h]} \frac{\partial v_p^*[h]}{\partial \omega[i, h]} \right) \\ \frac{\partial E}{\partial w[i, h]} &= - \sum_p ((y_p - \bar{y}_p) x[h] v_p[h] (1 - v_p[h]) u_p[i]) \end{aligned}$$

This changes for
different activation
functions

- Feed forward equations
 - $v_p^*[h] = \sum_i \omega[i, h] u_p[i];$
 - $v_p[h] = \frac{1}{1 + e^{-v_p^*[h]}};$
 - $y_p = b + \sum_h v_p[h] x[h]$
- $\frac{\partial v_p^*[h]}{\partial w(i, h)} = u_p[i]$
- $\frac{\partial v_p[h]}{\partial v_p^*[h]} = \frac{\partial (1 + e^{-v_p^*[h]})^{-1}}{\partial v_p^*[h]}$

$$= \frac{-1}{[1 + e^{-v_p^*[h]}]^2} e^{-v_p^*[h]} (-1)$$

$$= \frac{1 + e^{-v_p^*[h]} - 1}{[1 + e^{-v_p^*[h]}]^2} \text{ (add + 1, - 1 in numerator)}$$

$$= \frac{1 + e^{-v_p^*[h]}}{[1 + e^{-v_p^*[h]}]^2} - \frac{1}{[1 + e^{-v_p^*[h]}]^2}$$

$$= \frac{1}{1 + e^{-v_p^*[h]}} - \frac{1}{[1 + e^{-v_p^*[h]}]^2}$$

$$= v_p[h] (1 - v_p[h])$$

Backprop algorithm

1. Initialize

1. Initialize $b, \omega(\cdot, \cdot)$, and $x(\cdot)$ to random values
2. Set SSE_{old} to very large value.
3. Set $m = 0$ (iteration number)

2. Compute $v_p^*[h] = \sum_i w[i, h]u_p[i]$ for each p and each h .

3. Compute $v_p[h] = \frac{1}{1+e^{-v_p^*[h]}}$ for each p and each h

4. Compute $\bar{y}_p = b + \sum_h x[h]v_p[h]$ for each p (data samples)

5. Apply SD transformations ($E = \frac{1}{2} \sum_{p=1}^P (y_p - \bar{y}_p)^2$)

$$b_m \leftarrow b_{m-1} - \mu \frac{\partial(E)}{\partial b};$$

$$\frac{\partial E}{\partial b} = -\sum_p (y_p - \bar{y}_p)$$

6. $\omega_m[i, h] \leftarrow \omega_{m-1}[i, h] - \mu \frac{\partial(E)}{\partial \omega[i, h]};$

$$\frac{\partial E}{\partial \omega[i, h]} = -\sum_p ((y_p - \bar{y}_p)x[h]v_p[h](1 - v_p[h])u_p[i])$$

$$x_m[h] \leftarrow x_{m-1}[h] - \mu \frac{\partial(E)}{\partial x[h]};$$

$$\frac{\partial E}{\partial x[h]} = -\sum_p (y_p - \bar{y}_p)v_p[h]$$

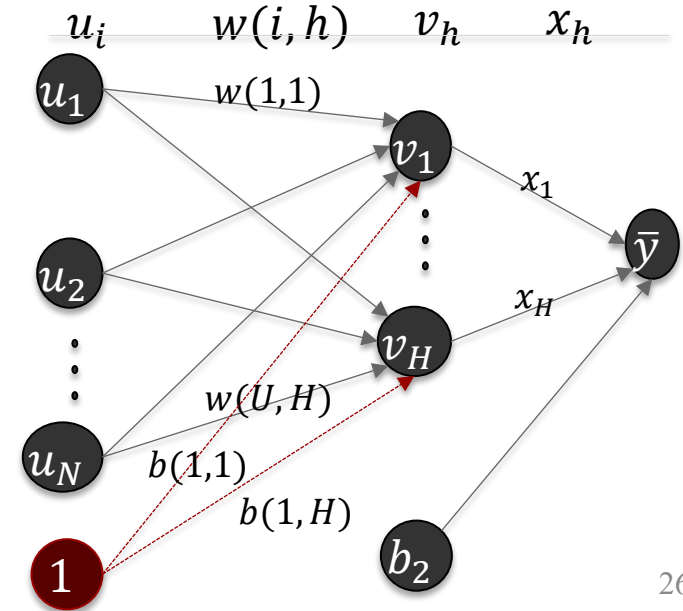
7. Set $m = m + 1$

1. Calculate $E_{new} = \sum_p (y_p - \bar{y}_p)^2$
2. Update μ
3. If $[E_{new} - E_{old}] < \text{tolerance} \rightarrow STOP$ Otherwise set $E_{old} = E_{new}$ got step 2.

- Do multiple iterations of the algorithm to start at different initial points
- Initialize $b, \omega[\cdot, \cdot]$, and $x[\cdot]$ to select from a random range; Try a few different ranges
- Try a few different options for learning rate ($\mu = \text{small constant}; \mu = \frac{A}{B+m}; \mu = \frac{1}{m};$)

This applies to sigmoid activation; change if using a different activation function

Backprop for below (addition of a bias node in input layer)?



Visualization of NN

- <https://cs.stanford.edu/people/karpathy/convnetjs/demo/regression.html>
- <https://cs.stanford.edu/people/karpathy/convnetjs/>
- <https://playground.tensorflow.org>

Further study

- Activation functions:
 - Sigmoid, Tanh, (for function approximation); vanishing gradient/exploding gradient challenges when layers increase
 - ReLU (overcomes vanishing gradient/exploding gradient problem)
 - Numerous other AF: <https://arxiv.org/abs/1811.03378>
- Neural network architectures:
 - <https://towardsdatascience.com/the-mostly-complete-chart-of-neural-networks-explained-3fb6f2367464>
- Optimizers (Backprop is typically the core, but estimations of gradient and types of learning rates may vary)
 - Adam, RMSProp, SGDm AdamW, Adadelata, Adamax, Adafactor, Nadam;
 - <https://keras.io/api/optimizers/>
 - <https://pytorch.org/docs/stable/optim.html>
 - Schmidt, et.al., ICML 2021; <https://github.com/SirRob1997/Crowded-Valley---Results>
- Deep learning: Sarker, I.H. Deep Learning: A Comprehensive Overview on Techniques, Taxonomy, Applications and Research Directions. *SN COMPUT. SCI.* **2**, 420 (2021).
<https://doi.org/10.1007/s42979-021-00>

Output layer

- Multiple nodes in output layer → when predicting multinomial functions
 - AF may then be used in output layer as well,
 - In RL, for policy gradient, we will use multiple nodes in output layer with softmax AF.
 - Lapan Ch3-Module1 code uses softmax
- Dropout in output layer :
 - It is a form of regularization in supervised learning, especially when data is sparse;
<https://arxiv.org/pdf/1207.0580.pdf>
 - It is not used that much in reinforcement learning, but there is some research in its use in policy gradient <https://doi.org/10.48550/arXiv.2202.11818>
 - Lapan Ch3-Module1 code uses dropout

Do you know?

- What we looked at is the simplest neural net architecture
- NN is a **single layer perceptron** developed in 1950's
- Deep learning- generally multi-layer NN are called DL
- **We looked at specific DL called feed forward,**
- Other types of deep learning (potentially a course in Fall 2025)
 - CNN (convolution neural network)
 - RNN (recurrent neural network)
 - Autoencoders
 - Attention networks
 - Self-attention mechanisms (Transformers)
- Single layer feedforward is very powerful for representation of most non-linear functions of any form – Universal approximators
- Other architectures were developed for image and text processing, processing of sequence data
 - Sequence based methods can be used for dynamic decision making (when we get to RL we will discuss deep learning)

Work colleague: I saw deep learning on your resume, I have a question for you..

New graduate:



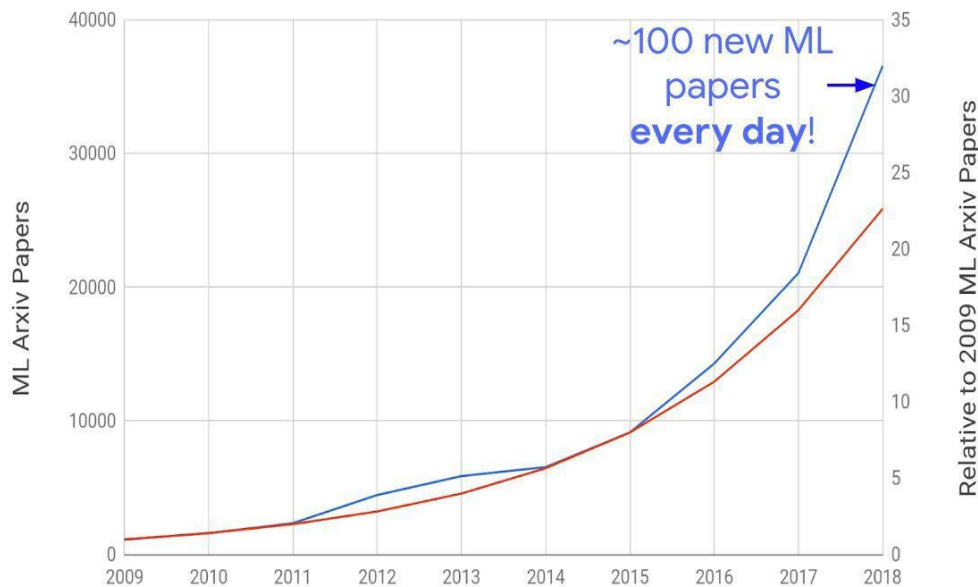
Work colleague: Great, we are debating on whether CNN or RNN for this what do you think?



Let's try to figure it out

Fast growing; hard to keep up with everything

— ML Arxiv Papers — Moore's Law growth rate (2x/2 years)



Student asks question



*Professor's response:
Just restates the question in different words*

If your mathematical foundation is strong, you can figure it out

Questions

- Have one hidden layer but increase number of nodes
- Vs.
- Have multiple hidden layers with fewer nodes in each
- The first does piecewise non-linear fits, while the second does non-linear transformations to a higher level;
 - So for a complex (highly non-linear function), number nodes needed to get same fit will be much higher than adding multiple layers.
 - The first would also be overfitting

Application areas of NN

- As a machine learning tool (**static** applications)– prediction, classification, function approximation
 - This application will follow typical machine learning rules for model training and accuracy; split by test and train data; evaluate predictive power using test scores etc.
 - How neural networks learn (insight into application of NN)
 - <https://www.youtube.com/watch?v=TkwXa7Cvfr8>
- As function approximation of environment (**dynamic** system representation)
 - e.g., recurrent neural network (RNN) can be trained to replace large-scale high-computation simulation models
 - Computational burden of training NN is high, but once trained computation burden is minimal
- As **function approximation for learning the policy function or value function in control optimization /intelligent systems** (**dynamic** decision making – we will see in deep RL)
 - This is the focus of this class
 - We will skip typical steps of ML training and testing here and look at overall model performance with RL

UMassAmherst
The Commonwealth's Flagship Campus