

# THEORETICAL FOUNDATIONS OF MACHINE LEARNING FOR ECONOMISTS: LECTURE 3

SONAN MEMON  
LECTURER, INSTITUTE OF BUSINESS  
ADMINISTRATION, KARACHI

14TH APRIL 2021

# OUTLINE OF LECTURE 3

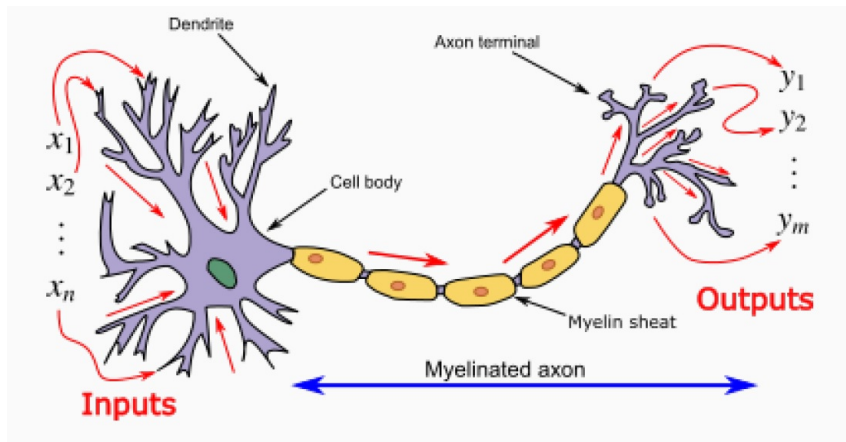
- i. Introduction to Deep Feed Forward Neural Networks
- ii. Gradient Based Learning
- iii. Output Units and Hidden Units
- iv. Network Architecture Design
- v. Backpropagation and Stochastic Gradient Descent
- vi. Regularization of DNN's
- vii. Applications in Economics
- viii. Implementation in R



## MOTIVATION

- i. Neural Networks are complicated function approximation methods, well suited for high dimensional and nonlinear settings.
- ii. Are *Networks* because there is a composition of many functions involved in the function approximation method.
- iii. Are *Deep* since the number of layers or the functions to be composed could be quite large.
- iv. Are *Neural* because of a loose analogy with neuroscience: hidden units in hidden layers, just like real neurons get inputs from many other units in last hidden layer, which is used to determine their activation.

# BIOLOGICAL MOTIVATION



## APPLICATIONS

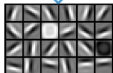
- i. AlphaGo.
- ii. Object recognition, Facial Recognition and Computer Vision: Convolutional Networks (form of deep, feed forward network).
- iii. Prediction of stock prices.
- iv. Recurrent Networks used for natural language processing (backward connections) such as translation, speech or digit recognition and syntactic parsing.
- v. Economic Applications: Solution of Nonlinear DSGE Models, High Dimensional Econometrics (Causal AI), High Dimensional Forecasting, Natural Language Processing and Computational Linguistics.

# ALPHAGO: THE GREATEST GO PLAYER OF ALL TIME?

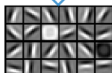


# OBJECT RECOGNITION EXAMPLES

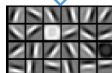
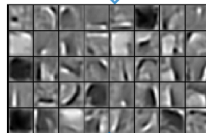
Faces



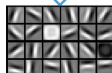
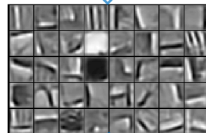
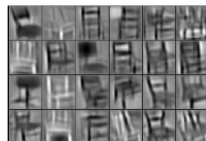
Cars



Elephants



Chairs





# INTRODUCTION TO DEEP NEURAL NETWORKS

- i. Choosing the cost function, form of hidden and output units and optimizer.
- ii. Network Design: Activation Functions, Network Architecture and Output Layers.
- iii. Calculating Gradients for Optimization: Back propagation and Stochastic Gradient Descent.
- iv. Regularization: Penalty Terms and Early Stopping.

# DEEP FEED FORWARD NETWORKS

- i. Deep Feed Forward Networks or MLP's are deep learning models in which information only flows forward and no *feedback*, unlike recurrent networks.
- ii. The goal is to learn unknown function  $y = f^*(\mathbf{x})$ .
- iii.  $y = f(\mathbf{x}) = f^n(f^{n-1}(f^{n-2} \dots (f^1(\mathbf{x}))))$
- iv. "Deep" here refers to the number of layers which are  $n$  in equation above.
- v. Width of network refers to the number of *hidden units* in one layer.
- vi. Hidden layers: only observe data on  $\mathbf{x}$  and  $y$ , not on intermediate layers, which have to be determined by the learning algorithm.

# FLOW REPRESENTATION

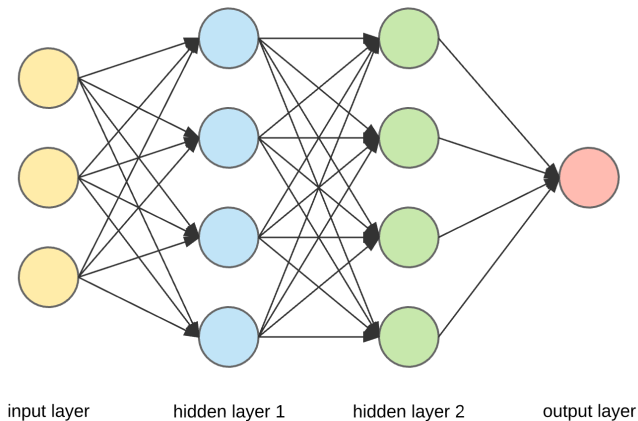


Figure: Flow Representation

# DEEP FEED FORWARD NETWORKS

- i. Simple linear model of form  $f = \mathbf{w}^T \mathbf{x} + b$  cannot model interactions between inputs and has limited model capacity.
- ii. We apply nonlinear transformation  $\phi(\mathbf{x})$  before applying linear model to the output of  $\phi$  to overcome limitations of linear models:  $y = f(\mathbf{x}; \theta, \mathbf{w}) = \phi(\mathbf{x}; \theta)^T \mathbf{w}$ .
- iii. How to choose  $\phi$ ?:
  - ▶ Generic  $\phi$
  - ▶ Manual engineering through specialized human effort.
  - ▶ Learn  $\phi$  using deep learning.
- iv. Using deep learning, we optimally learn parameters  $\theta$  to pin down  $\phi$  and this combines the benefits of approaches 1 and 2 above.



## EXAMPLE: LEARNING XOR

- i. We want our network to perform exactly correctly on the four points  $\mathbb{X} = \{[0, 0]^T, [0.1]^T, [1, 0]^T, [1, 1]^T\}$  in training set and no statistical generalization is needed here.
- ii. Evaluated on our whole training set, the MSE loss function is:  

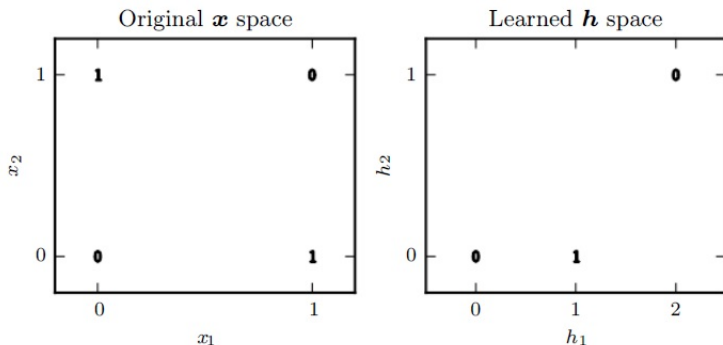
$$L(\theta) = \frac{1}{4} \sum_{\mathbf{x} \in \mathbb{X}} (f^*(\mathbf{x}) - f(\mathbf{x}; \theta))^2.$$
- iii. The linear model cannot capture the nonlinearity of this function.

## EXAMPLE: LEARNING XOR

- i. We use feed forward network with one hidden layer by defining  $f^1(\mathbf{x}) = \mathbf{W}^T \mathbf{x} + \mathbf{c}$  and  $f^2(\mathbf{h}) = \mathbf{w}^T \mathbf{h} + b$ ,  $\mathbf{h} = g(\mathbf{W}^T \mathbf{x} + \mathbf{c})$ .
- ii. If  $g = \text{ReLU}$  activation function,  

$$f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^T \max\{0, \mathbf{W}^T \mathbf{x} + \mathbf{c}\} + b.$$
- iii. The hidden layer does feature extraction and a linear model is applied to transformed inputs after applying activation function  $g$ .
- iv. The initial untransformed space cannot be learned by the linear function but the transformed feature space can be learned by linear function.

# ORIGINAL SPACE VERSUS LEARNED SPACE





## EXAMPLE: LEARNING XOR

- i.  $\mathbf{W}$ ,  $\mathbf{c}$ ,  $b$ ,  $\mathbf{w}$  can be chosen to exactly find the function  $f^*$  in this case.
- ii.  $\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ ,  $\mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$ ,  $\mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$ ,  $b = 0$ .
- iii. With above parameters, the neural network learns every example in the batch perfectly, with zero error.
- iv. In a real situation, there might be billions of model parameters and billions of training examples, implying that a gradient-based optimization algorithm can find parameters that produce very little error but not zero error.

# GRADIENT BASED LEARNING

- i. Neural networks are usually trained by using iterative, gradient-based optimizers that merely drive the cost function to a very low value, rather than zero value.
- ii. Backpropagation and stochastic gradient descent is used for computational efficiency.
- iii. Stochastic gradient descent applied to non-convex loss functions has no convergence guarantee, and is sensitive to the values of the initial parameters.
- iv. Computing the gradient for neural network models is more complicated but fundamentally similar to gradient based learning of linear models.

# LOSS FUNCTION CHOICE

- i. Parametric model defines conditional distribution  $p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$  and we use maximum likelihood for loss function, which is simply the negative log-likelihood (cross-entropy between data and model):  $L(\boldsymbol{\theta}) = -\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{data}} \log p_{model}(\mathbf{y}|\mathbf{x})$
- ii. The loss function used in practice will often add a regularization term.
- iii. Functions that saturate make the gradient vanish and inhibit learning. The negative log-likelihood helps to avoid this problem because log in the loss function undoes the exp of some output units such as softmax and sigmoid.

# OUTPUT UNITS

- i. The choice of cost function is tightly coupled with the choice of output unit since how we represent the output determines the form of the likelihood function.
- ii. The role of the output layer is to apply additional transformation on the features in hidden layers:  $\mathbf{h} = f(\mathbf{x}; \boldsymbol{\theta})$ , in line with nature of task.
- iii. Linear units, sigmoid units and softmax units are some commonly used output units.

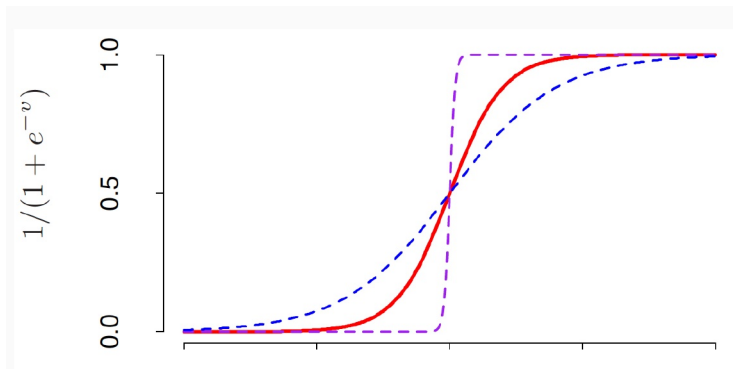
# LINEAR OUTPUT UNITS

- i. Simple output units are affine transformation of the hidden units. For instance, given features  $h$ , linear output unit will be  $\hat{y} = \mathbf{W}^T \mathbf{h} + \mathbf{b}$ .
- ii. Because linear units do not saturate, they pose little difficulty for gradient based optimization algorithms and may be used with a wide variety of optimization algorithms.
- iii. Used for regression problems.

## SIGMOID OUTPUT UNITS

- i. Many tasks require predicting the value of a binary variable  $y$  such as classification problems with two classes.
- ii. A Bernoulli distribution is defined by the conditional probability of success, which is what the neural net predicts:  $\mathbb{P}(y = 1|\mathbf{x})$ .
- iii. A common solution in this case is to use sigmoid output units, applied on affine transformation of  $\mathbf{h}$ :  $y = \sigma(\mathbf{w}^T \mathbf{h} + b)$ , where  $\sigma$  is logistic sigmoid function.

# SIGMOIDAL FUNCTION



## SOFTMAX OUTPUT UNITS

- i. When we wish to represent a probability distribution over a discrete variable with  $n$  possible values, we may use the softmax function.
- ii. We need to produce  $\hat{\mathbf{y}} \in \mathbb{R}^n$  such that  $\hat{y}_i = \mathbb{P}(y = i | \mathbf{x})$ , where  $\sum_{i=1}^n \hat{y}_i = 1$ .
- iii.  $z_i = \log \tilde{P}(y = i | \mathbf{x})$  is the log of un-normalized probability.
- iv. First apply affine transformation to get  $\mathbf{z} = \mathbf{W}^T \mathbf{h} + \mathbf{b}$ , then apply softmax function element-wise to get:  

$$\text{softmax}(z_i) = \frac{\exp\{z_i\}}{\sum_{j=1}^n \exp\{z_j\}}$$
 (exponentiation plus normalization).



# HIDDEN UNITS

- i. Hidden units refer to the activation functions used in hidden layers.
- ii. In the 1990's, logistic/sigmoid functions were default choices and hyperbolic tangents took over in late 1990's to 2000's.
- iii. Rectified Linear Units (ReLU) are a default choice of hidden units in the modern ML literature since it solves the vanishing gradient problem etc.
- iv.  $\mathbf{h} = g(\mathbf{W}^T \mathbf{x} + \mathbf{b})$ , where  $g(z) = \max\{0, z\}$  is the ReLU activation function.
- v.  $g(z) = \log(1 + e^z)$  is the softplus activation function. It has been used in recent econ literature.

# SIGMOID AND HYPERBOLIC TANGENT

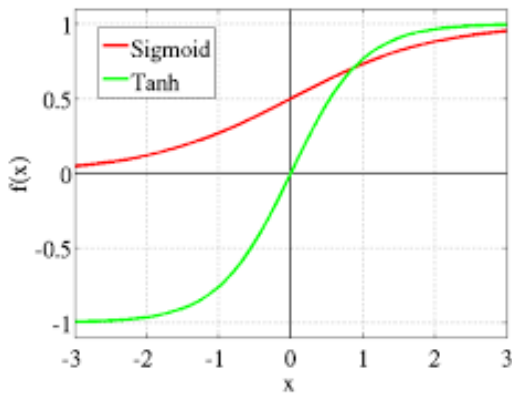
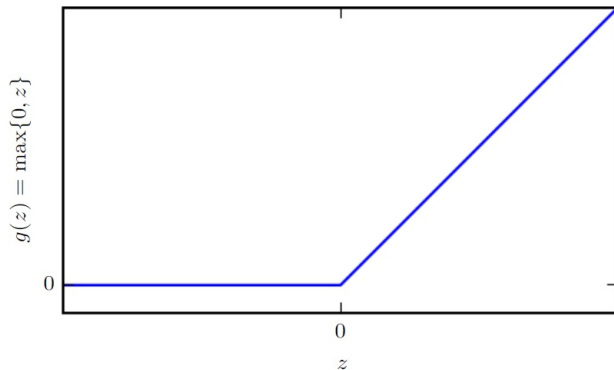
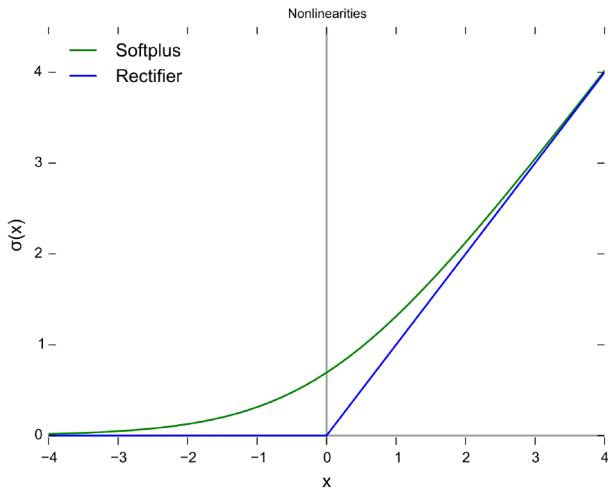


Figure: Sigmoid and Tanh

# RELU ACTIVATION FUNCTION



# ACTIVATION FUNCTIONS: RELU AND SOFTPLUS



# NETWORK ARCHITECTURE DESIGN

- i. The *depth*, *width* and structure of connections between hidden units across layers are key design features of the architecture of a neural network.
- ii. Depth refers to number of layers and width refers to number of hidden units per layer.
- iii. For instance,  $\mathbf{h}^1 = g^1 \left( (\mathbf{W}^1)^T \mathbf{x} + \mathbf{b}^1 \right)$  may be first layer and  $\mathbf{h}^2 = g^2 \left( (\mathbf{W}^2)^T \mathbf{x} + \mathbf{b}^2 \right)$  may be second layer etc.
- iv. Deeper networks can use far fewer units per layer and far fewer parameters to fit data and generalize to the test set better, but are also often harder to optimize.
- v. Even one layer with sufficient parameters and units can fit the data in theory.

# UNIVERSAL APPROXIMATION THEOREM

Theorem (Universal Approx Theorem Hornik et al 1989)

A feed forward network with a linear output layer and at least one hidden layer with any “squashing” activation function such as logistic sigmoid can approximate any Borel measurable function from one finite-dimensional space to another with any desired non-zero amount of error, provided that the network is given enough hidden units.

# UNIVERSAL APPROXIMATION THEOREM

- i. A feedforward network with a single layer or few layers is sufficient to represent any function in theory but in worst cases, the number of units required in the one or few layers may be prohibitively large.
- ii. There is no universally superior optimization algorithm to be applied on a particular training data set, which means that we may fail to learn parameters corresponding to the true function.
- iii. Failures may also occur due to over fitting or failure to generalize correctly.

## WHY DEEPER MODELS?

- i. In many circumstances, using deeper models can reduce the number of units required to represent the desired function and can reduce the amount of generalization error.
- ii. From a Bayesian perspective, deeper models may reflect our priors that there are many but simpler underlying factors of variation which when composed together constitute the complexity of what has to be learned.



## HOW TO CONNECT TWO LAYERS?

- i. When we apply affine transformations of form  $\mathbf{W}^T \mathbf{x} + \mathbf{b}$  in networks, the matrix  $\mathbf{W}$  will in general connect every input unit to every output unit.
- ii. Many specialized networks have fewer connections, so that each unit in the input layer is connected to only a small subset of units in the output layer.
- iii. This sparsity can reduce the number of parameters and the amount of computation required to evaluate the network but specific pattern of sparsity that is efficient is highly problem-dependent.

## BACKPROPAGATION (BAKCPROP)

- i. Numerically evaluating the gradient for optimization can be computationally expensive and backprop algorithm does so using a simple and inexpensive procedure.
- ii. In a feed forward network information flows forward from inputs  $\mathbf{x}$ , through hidden layers, ultimately producing loss  $L(\theta)$ .
- iii. backprop allows the information from the loss to then flow backwards through the network, in order to compute the gradient.
- iv. backprop is not an learning algorithm but an efficient method for computing derivative, which is an intermediate step in the full learning algorithm such as stochastic gradient descent.

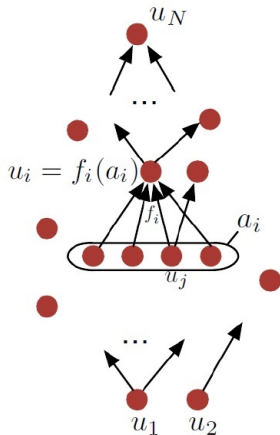
## BACKPROPAGATION: CHAIN RULE

- ii. A straightforward computation of chain rule introduces many subexpressions which have to be recomputed several times and these wasteful computations may even grow exponentially.
- iii. backprop performs the chain rule operation for each “edge” of the graph and visits each edge from node  $u^i$  to  $u^j$  *exactly once*.
- iv. backprop stores intermediate derivatives rather than recomputing them, which when combined with its recursive structure improves efficiency.

# BACKPROPAGATION: COMPUTATIONAL GRAPH

- i. Consider computational graph in which we want to compute  $\frac{\partial u^n}{\partial u^i}$ : gradient of  $u^n$  (e.g loss associated with mini-batch) with respect to  $M$  input nodes  $u^i$  where  $i \in \{1, 2, \dots, M\}$ .
- ii. We will assume that the nodes of the graph have been ordered in such a way that we can compute their output one after the other and  $u^i = f(A^i)$ , where  $A^i$  are set of parent nodes of node  $u^i$ .
- iii. Chain rule implies that  $\frac{\partial u^n}{\partial u^j} = \sum_{i:j \in Pa(u^i)} \frac{\partial u^n}{\partial u^i} \frac{\partial u^i}{\partial u^j}$ .
- iv. The forward computation step shown next has already stored  $\frac{\partial u^n}{\partial u^i}$  and the computation for each edge corresponds to computing a partial derivative  $\frac{\partial u^i}{\partial u^j}$ , performing one multiplication and one addition only.

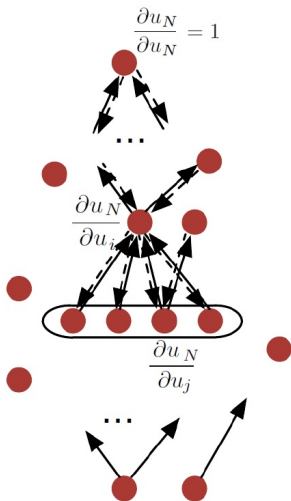
# RECURSIVE FORWARD COMPUTATION



# RECURSIVE FORWARD COMPUTATION

```
for  $i = 1, \dots, M$  do  
   $u_i = x_i$   
end for  
for  $i = M + 1, \dots, N$  do  
   $a_i = (u_j)_{j \in \text{parents}(i)}$   
   $u_i = f_i(a_i)$   
end for  
return  $u_N$ 
```

# RECURSIVE BACKWARD COMPUTATION



## RECURSIVE BACKWARD COMPUTATION

```

 $\frac{\partial u_N}{\partial u_N} = 1$ 
for  $j = N - 1, \dots, 1$  do
   $\frac{\partial u_N}{\partial u_j} = \sum_{i: j \in \text{parents}(i)} \frac{\partial u_N}{\partial u_i} \frac{\partial u_i}{\partial u_j}$ 
end for
return  $\left( \frac{\partial u_N}{\partial u_i} \right)_{i=1}^M$ 

```



## REGULARIZATION FOR DEEP NEURAL NETWORKS

- i. As in other machine learning algorithms, improving predictive performance requires regularization, which can be done by introducing  $L^1$  (LASSO) or  $L^2$  (Ridge Regression) penalty.
- ii. However, tuning using cross-validation is computationally demanding with deep nets.
- iii. Alternative regularization method, commonly used is **early stopping**:
  - i. Split the data into training and validation samples.
  - ii. Run gradient based optimization on training sample and check prediction loss in validation sample at each iteration.
  - iii. Stop when prediction loss starts increasing.

## OPTIMIZATION: GRADIENT DESCENT

- i. The algorithm involves updating the model parameters  $\theta$  which would include the weights and biases associated with each layer with a small step in the direction of the gradient of the loss function.
- ii. With training data  $(\mathbf{x}^t, \mathbf{y}^t)$ ,  

$$\theta = \theta + \epsilon \nabla_{\theta} \sum_t L(f(\mathbf{x}^t; \theta), \mathbf{y}^t; \theta).$$
- iii. Choice of learning rate is crucial since a too low learning rate slows down learning and too high value can cause divergence or oscillation around the minima.

# GRADIENT DESCENT

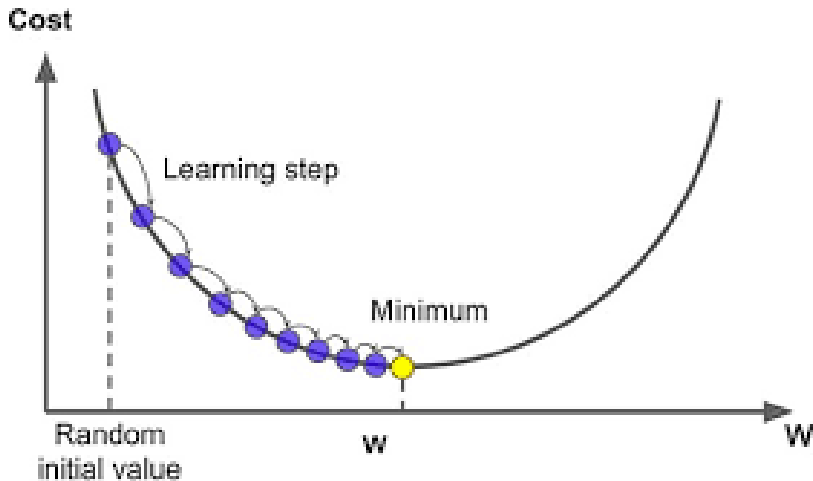


Figure: Learning via GD

## WHY STOCHASTIC GRADIENT DESCENT (SGD)?

- i. SGD uses just a random subsample of the data for gradient updating.
- ii. The key observation that motivates SGD is that in an (i.i.d.) sampling context, further observations are increasingly redundant.
- iii. Formally, the standard error of a gradient estimate based on  $b$  observations is of order  $\frac{1}{\sqrt{b}}$  but computation time is of order  $b$ .
- iv. For a large dataset, it would take forever to just calculate one gradient, and do one updating step. During the same time, SGD might have made many steps and come considerably closer to the truth (see Bottou et al. (2018))

# STOCHASTIC GRADIENT DESCENT ALGORITHM

Note:  $\eta$  is learning rate

**while** stopping criterion not met **do**

sample **minibatch** of  $m$  examples from the training set

$\{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^m\}$ .

set  $\mathbf{g} = 0$

**for**  $t = 1 : m$  **do**

$\mathbf{g} = \mathbf{g} + \nabla_{\theta} L(f(\mathbf{x}^t; \theta), \mathbf{y}^t; \theta)$

**end for**

Apply update:  $\theta = \theta - \eta \mathbf{g}$

**end while**

## SGD WITH MOMENTUM

- i. Momentum is a method that helps accelerate SGD in the relevant direction and dampens oscillations of SGD.
- ii. The momentum term is usually set to 0.9 or a similar value.
- iii. The momentum term increases for dimensions whose gradients point in the same directions and reduces updates for dimensions whose gradients change directions. As a result, we gain faster convergence and reduced oscillation.

# STOCHASTIC GRADIENT DESCENT WITH MOMENTUM

Note:  $\eta$  is learning rate, momentum parameter is  $\alpha$   
Initial parameter  $\theta$  and initial velocity  $\mathbf{v}$  must be specified.

**while** stopping criterion not met **do**

sample **minibatch** of  $m$  examples from the training set  $\{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^m\}$ .

set  $g = 0$

**for**  $t = 1 : m$  **do**
$$\mathbf{g} = \mathbf{g} + \nabla_{\theta} L(f(\mathbf{x}^t; \theta), \mathbf{y}^t; \theta)$$
**end for**

velocity update:  $\mathbf{v} = \alpha \mathbf{v} - \eta \mathbf{g}$

apply update:  $\theta = \theta + \mathbf{v}$

end while

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ 🔍 ↺



# ECONOMIC APPLICATION 1: COMPUTATION OF NONLINEAR DSGE MODELS

- i. Households are subject to uninsurable idiosyncratic risk and can only save in the risk-free bond to insure against idiosyncratic shocks. Meanwhile, financial experts absorb all the capital return risk.
- ii. The interaction between the supply of bonds by the financial sector and the precautionary demand for bonds by households produces significant endogenous aggregate risk and an endogenous regime-switching for real variables including leverage across the two stochastic steady states of model (SSS).

# ECONOMIC APPLICATION 1: COMPUTATION OF NONLINEAR DSGE MODELS

- i. Endogenous regime switching generates bimodal and skewed ergodic distributions of aggregate variables, time-variation in volatility and skewness, and supercycles of borrowing and deleveraging.
- ii. In the high (low) leverage SSS, endogenous aggregate risk is high (low) and after a negative aggregate shock, the economy suffers deep and protracted (mild) recessions with persistently low wages (quick recovery) since financial sector loses a lot of its wealth.

# ECONOMIC APPLICATION 1: COMPUTATION OF NONLINEAR DSGE MODELS

- i. The wealth distribution is an infinite-dimensional object so standard DP techniques cannot be employed.
- ii. The consumption decision rule of households is sharply nonlinear with respect to the aggregate state variables (equity and debt).
- iii. The global, nonlinear solution method used is ML: DNN to obtain a flexible, nonlinear approximation of the perceived law of motion of the cross-sectional distribution of assets with finite moments.
- iv. Neural networks can capture a nonlinear PLM without having to specify, ex-ante, any concrete structure for it.

# ECONOMIC APPLICATION 1: COMPUTATION OF NONLINEAR DSGE MODELS

- i. By solving the model (read paper for derivation), they derive that  $dB_t = (w_t + r_t B_t - C_t)d_t$  is the actual law of motion of debt in the model.
- ii. Given their model, it is crucial to find the households' aggregate consumption  $C_t$  to compute the equilibrium of the economy given the structural parameters.
- iii. Follow [Krusell and Smith \(1998\)](#) and assume that households only use a finite set of  $n$  moments of the cross-sectional distribution of assets instead of complete distribution for expectation formation.
- iv. The PLM with  $n = 1$  is  $dB_t = h(B_t, N_t)d_t$  where 
$$h(B_t, N_t) = \frac{\mathbb{E}[dB_t|B_t, N_t]}{dt}.$$

# NONLINEAR DSGE MODELS

- i. The functional form  $h(B_t, N_t)$  is approximated by applying NN's to simulated data from the aggregate ergodic distribution.
- ii. NN's help approximate the PLM arbitrarily well, allowing for complex nonlinearities in the law of motion of  $B_t$  and  $N_t$  in the model.
- iii. NN's also allow for good extrapolation outside their training areas of the state space unlike other methods using Chebyshev polynomials.

# ECONOMIC APPLICATION 1: COMPUTATION OF NONLINEAR DSGE MODELS

- i. A neural network with a single, hidden layer is used:  
$$h(\mathbf{s}; \boldsymbol{\theta}) = \theta_0^2 + \sum_{q=1}^Q \theta_q^2 \phi \left( \theta_0^1 + \sum_{i=1}^N \theta_i^1 s^i \right).$$
- ii. ReLU is common choice in ML literature but they use  $\phi(x) = \log(1 + e^x)$  (softplus) to avoid the kink in ReLU.
- iii. Given a choice of  $Q$  (number of hidden units), input vector of states  $\mathbf{s}$  and choice of number of layers and loss function, the  $\theta$  parameters can be estimated, which help approximate the PLM.
- iv. They use batch gradient (not SGD) algorithm in which loss gradient is computed using backprop to update  $\theta$  parameters.

## ECONOMIC APPLICATION 2: COUNTERFACTUAL PREDICTION WITH DEEP IV NETWORKS

- i. Decision makers look to the data to model counterfactual questions and ML will do a poor job of predicting the many potential futures associated with each policy option.
- ii. For example, optimal pricing requires predicting sales under *changes* to prices.
- iii. In order to accurately answer such counterfactual questions it is necessary to model the structural relationship between policy and outcome variables.
- iv. This paper uses the concept of IV's to construct systems of ML tasks that can be applied in causal inference on large and unstructured data.

## ECONOMIC APPLICATION 2: COUNTERFACTUAL PREDICTION WITH DEEP IV NETWORKS

- i. Set up IV system of machine learning tasks that can be targeted with deep learning and allow us to make counterfactual claims and perform causal inference.
- ii. Two ML stages: a first stage that models the conditional distribution for treatment given the instruments and covariates, and a second stage which targets a loss function involving integration over the conditional treatment distribution from the first stage.
- iii. Both stages use deep neural nets trained via stochastic gradient descent.
- iv. Out of sample causal validation as regularization also performed.



## ECONOMIC APPLICATION 2: COUNTERFACTUAL PREDICTION WITH DEEP IV NETWORKS

- i.  $y = g(p, x) + e$ , where  $\mathbb{E}[pe|x] \neq 0$ .
- ii.  $h(p, x) = g(p, x) + \mathbb{E}[e|x]$
- iii. To evaluate policy options, say  $p_0$  and  $p_1$ , we can look at the difference in mean outcomes:  

$$h(p_1; x) - h(p_0; x) = g(p_1; x) - g(p_0; x).$$
- iv. In standard, unstructured ML, the prediction model is trained to fit  $\mathbb{E}[y|p, x] = g(p, x) + \mathbb{E}[e|p, x] \neq h(p, x)$  since  $\mathbb{E}[e|p, x] \neq \mathbb{E}[e|x]$  so that ML will lead to biased counterfactual analysis.
- v. IV's allow us to get unbiased estimate  $\widehat{h(p, x)}$ .

## ECONOMIC APPLICATION 2: COUNTERFACTUAL PREDICTION WITH DEEP IV NETWORKS

- i. This paper avoids explicit linearization used by standard IV procedures such as 2SLS by instead using DNN's, optimized via SGD to learn  $F$  and  $h$ .
- ii. In first stage, we learn  $F(p|x, z)$  by choosing appropriate distribution, parametrized by  $\hat{F} = F_{\theta}(p|x, z)$ .
- iii. Estimation proceeds by maximum likelihood via stochastic gradient descent on the implied negative log likelihood.
- iv. In the second stage, conditional on estimate of  $F$ , the counterfactual function  $h$  is approximated by a DNN, parametrized as  $h_{\theta}$ , based on minimization of some loss function.

## ECONOMIC APPLICATION 2: COUNTERFACTUAL PREDICTION WITH DEEP IV NETWORKS

- i. Each stage is evaluated in turn, with second stage validation using the best-possible network as selected in the first stage.
- ii. In the first stage, cross-validation avoid overfitting and guards against weak instruments.
- iii. In the second stage, validation will lead to best-possible performance on the objective  $\mathbb{E}[y|x, z]$  under the constraint imposed by  $F_\theta$  from first stage.
- iv. Dubbed as Causal AI or Artificial Economic Intelligence.

## IMPLEMENTATION IN R

- i. [playground.tensorflow.org](http://playground.tensorflow.org) is a website where you can tweak and visualize neural networks.
- ii. Many packages exist in R for training deep, neural networks such as `mxnet`, `H2O`, `keras` etc.
- iii. `mx.` functions from `mxnet` package.
- iv. See <https://www.datacamp.com/community/tutorials/keras-r-deep-learning> for an interesting application using image data.

## PSEUDO CODE IN R

- i. library(mxnet)
- ii. Network Structure:
  - ▶ `data <- mx.symbol.Variable("data")`
  - ▶ 1st hidden layer with 128 neurons:  
`f1 <- mx.symbol.FullyConnected(data, name = "f1",  
numhidden = 128)`
  - ▶ Set the activation which takes output of  $f_1$ :  
`act1 <- mx.symbol.Activation(f1, name = "relu1", acttype =  
"relu")`
  - ▶ second layer takes output from  $act_1$  and then define second  
activation function and so on:  
`f2 <- mx.symbol.FullyConnected(act1, name = f2, numhidden =  
64)`
  - ▶ ...
  - ▶ finally output layer with 10 neurons:  
`f3 <- mx.symbol.FullyConnected(act2, name = "f3",  
numhidden = 10)`

## PSEUDO CODE IN R

- i. softmax output unit for multinomial output case:  
`softmax <- mx.symbol.SoftmaxOutput( $f_3$ , name = "sm")`
- ii. assign cpu to mxnet  
`devices <- mx.cpu()`
- iii. set seed to control the random process in mxnet  
`mx.set.seed(0)`

## PSEUDO CODE IN R

### Train the Neural Network:

- i. `model <- mx.model.FeedForward.create(softmax, X =  $train_x$ , y =  $train_y$ , ctx = devices, num.round = 10, array.batch.size = 100, learning.rate = 0.07, momentum = 0.9, eval.metric = mx.metric.accuracy, initializer = mx.init.uniform(0.07), epoch.end.callback = mx.callback.log.train.metric(100))`
- ii. num.round is number of iterations for training, batch size is batch size for SGD. Learning and momentum parameters are also for SGD. “initializer” is initialization scheme for parameters.
- iii. eval.metric is the evaluation function applied on the results to monitor performance on the validation set. “accuracy” calculates the mean accuracy per sample for softmax in all dimensions.
- iv. Make a prediction: `preds <- predict(model,  $test_x$ )`

## PSEUDO CODE IN R

- iv. callback function saves a checkpoint to files during each period iteration. “log.train.metric” logs a training metric each period.



**Thank you**