

# Lecture 1 Introduction

Shiwei Lan<sup>1</sup>

<sup>1</sup>School of Mathematical and Statistical Sciences  
Arizona State University

STP598 Advanced Deep Learning Models  
Spring 2026

Overview

S.Lan

Deep Learning:  
Opportunities  
and Challenges

Course Overview

DNN Review

Computing

## Advanced Deep Learning Models

- Instructor: Shiwei Lan <[slan@asu.edu](mailto:slan@asu.edu)>
  - Office: WXMLR 544
  - Office hours: TuTh 12 - 1:15p @ <https://asu.zoom.us/j/8055899886>
- Teaching Assistant:
  - Ping-Han Huang <[phuang37@asu.edu](mailto:phuang37@asu.edu)>
  - Zoom: <https://asu.zoom.us/j/87568448130?jst=2>
  - Office hours: TBA

Overview

S.Lan

Deep Learning:  
Opportunities  
and Challenges

Course Overview

DNN Review

Computing

- Course Schedule, Lecture Notes, etc.  
<https://slan-teaching.github.io/STP598adl/>
- Discussion and Questions  
Canvas Discussions [STP 598adl](#)
- Homework and Grades on [Canvas](#)  
<https://canvas.asu.edu/courses/248838>
- Coding Assignments on [Nbgrader](#)  
[https://agnesi.la.asu.edu/services/stp598\\_28057](https://agnesi.la.asu.edu/services/stp598_28057)

Overview

S.Lan

Deep Learning:  
Opportunities  
and Challenges

Course Overview

DNN Review

Computing

1 Deep Learning: Opportunities and Challenges

2 Course Overview

3 DNN Review

4 Computing

# What is Deep Learning?

Overview

S.Lan

Deep Learning:  
Opportunities  
and Challenges

Course Overview

DNN Review

Computing

**Deep learning** is a branch of machine learning that uses algorithms called **neural networks**—inspired by the human brain—to learn patterns from large amounts of data.



- A neural network is made of layers of connected “neurons” (mathematical functions).
- Deep learning means the network has many layers (input → multiple hidden layers → output).
- Each layer learns increasingly complex features:
  - In images: edges → shapes → objects
  - In language: letters → words → meaning
- The model learns by adjusting millions (or billions) of parameters to minimize errors.

**Deep learning** automatically *learns useful features from raw data*, instead of relying on humans to hand-engineer them.

- ① It handles complex, high-dimensional data
  - It works extremely well with images and video, Speech and audio, natural language (text, chatGPT), and sensor data (self-driving, robotics).
- ② It reduces manual feature engineering
  - Before: experts manually decide what features matter; After: the model learns the features by itself.
- ③ It improves accuracy at scale
  - When you have a lot of data and computing power, deep learning usually outperforms other methods. (Scaling law).
- ④ It enables modern AI applications
  - Image recognition (face unlock, medical imaging), voice assistants (Google, Alexa), FSD (Tesla), LLM (ChatGPT, Gemini, Claude).

## ① Vanishing and Exploding Gradients

- As networks became deeper, gradients during backpropagation would either vanish or explode. Solution: activation functions, initialization schemes, normalization, residual learning.

## ② Overfitting and Poor Generalization

- Large networks memorize training data and perform poorly on unseen data. Solution: regularization (dropout), data augmentation, pre-training+fine-tuning (transfer-learning), self-supervised learning.

## ③ Data Hunger and Label Dependence

- Self-supervised and contrastive learning (BERT-style masked modeling), foundation models.

## ④ Limited Long-Range Dependency Modeling

- CNNs, RNNs or LSTMs struggle with long sequences and global context. Solution: transformers and attention mechanisms, sparse and linear attention variants.



## 5 Training Instability and Optimization Difficulty

- Training is sensitive to learning rates, initialization, hyperparameters. Solution: advanced optimizers, learning rate schedules, and gradient clipping.

## 6 Computational Inefficiency

- Hardware acceleration (GPUs, TPUs, NPUs), parallelism, efficient architectures (MoE).

## 7 Lack of Interpretability

- DNNs are "black boxes" and hard to understand or trust. Solution: attention visualization, mechanistic interpretability, and physics informed learning.

## 8 Task-Specific, Narrow Models

- Trained for a single task, hard to reuse. Solution: foundation models, multitask and multimodal learning, and prompt-based adaptation.

Overview

S.Lan

Deep Learning:  
Opportunities  
and Challenges

Course Overview

DNN Review

Computing

- ① Generative AI (e.g., GPT, DALL-E, Stable Diffusion)
  - [generative AI: Generative artificial intelligence \(Wikipedia\)](#)
- ② ChatGPT and Advanced Language Models
  - [AI Breakthrough Timeline](#)
- ③ AlphaGo and Game-Playing AI
  - [70 Years of AI: Key Milestones](#)
- ④ AlphaFold — Protein Structure Prediction
  - [AlphaFold Changed Science](#)
- ⑤ Advances in Multimodal and Reasoning AI
  - [AI Breakthrough Timeline](#)

Overview

S.Lan

Deep Learning:  
Opportunities  
and Challenges

Course Overview

DNN Review

Computing

1 Deep Learning: Opportunities and Challenges

2 Course Overview

3 DNN Review

4 Computing

- **Residual Networks (ResNet)**: enabling training of very deep NNs
- **U-Net**: fundamental CNN architecture for image segmentation
- **General Adversarial Networks (GAN)**: introducing adversarial training to generate realistic images
- **Transformer**: introducing self-attention mechanism to enable parallelism and scaling, backbone of modern foundation models
- **Bidirectional Encoder Representations from Transformers (BERT)**: shifting NLP from task-specific models to more general
- **Diffusion Models**: dominant generative models
- **Graph Neural Networks (GNN)**: network for graph-structured data
- **Large language models (LLMs)**: modern AI applications

Overview

S.Lan

Deep Learning:  
Opportunities  
and Challenges

Course Overview

DNN Review

Computing

1 Deep Learning: Opportunities and Challenges

2 Course Overview

3 DNN Review

4 Computing

Overview

S.Lan

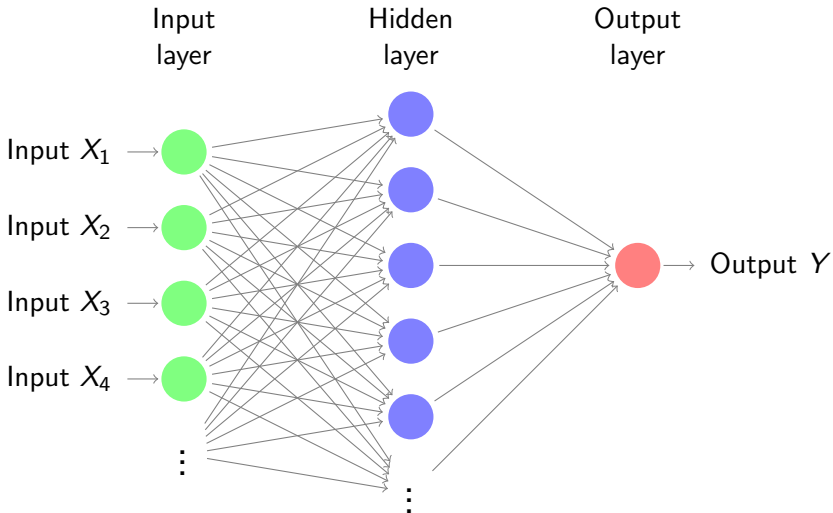
Deep Learning:  
Opportunities  
and Challenges

Course Overview

DNN Review

Computing

- Neural Networks were first developed as models for the human brain, where we have many units (**neurons**) that simultaneously process signals to give a joint decision.
- The neurons fire when the total signal passed to that unit exceeds a certain threshold.
- The collective signal from all neurons tells you whether its a dog or a cat.



- Given a training set  $\{x_i, y_i\}_{i=1}^n$ ,
  - For regression:  $y_i \in \mathbb{R}^K$  is a  $K$  dimensional continuous outcome
  - For classification:  $y_i \in \{1, 2, \dots, K\}$
- The goal is still to model the relationship

$$E(Y|X) = f(X)$$

- Instead of modeling the probabilities directly using  $X$ , we build  $M$  hidden neurons as a hidden layer between  $X$  and  $Y$ :

$$\begin{aligned} Z &= (1, Z_1, Z_2, \dots, Z_M) \\ &= (1, \sigma(X^T \alpha_1), \sigma(X^T \alpha_2), \dots, \sigma(X^T \alpha_M)) \end{aligned}$$



- $\sigma(\cdot)$  is an **activation function**. Some examples?
- We model  $Y$  using the hidden layer variables  $Z$  through some **link function**  $g(\cdot)$

$$X \xrightarrow{\sigma(\cdot)} Z \xrightarrow{g(\cdot)} Y$$

- In **classification problems** ( $K$  class), we can use logit link  $g_k$  to model the probability of  $Y = k$ , for  $k = 1, \dots, K$ :

$$g_k(Z) = \frac{\exp(Z^T \beta_k)}{\sum_{l=1}^K \exp(Z^T \beta_l)}$$

- In **regression problems** (could be multidimensional), we can simply use a linear function to model the  $k$ th entry of  $Y$ :

$$g_k(Z) = Z^T \beta_k$$

- The multidimensional function  $\mathbf{f}(x)$  can be represented as a convoluted way of mapping  $x \in \mathbb{R}^p$  to  $y \in \mathbb{R}^K$

$$\mathbf{f}(x) = \mathbf{g} \circ \sigma(x)$$

- The notations  $\mathbf{g}$  and  $\sigma$  here are multidimensional.
- The parameters involved are:  $\alpha_1, \dots, \alpha_M$ , and  $\beta_1, \dots, \beta_K$ .

- The **activation function**  $\sigma(\cdot)$  takes a linear combination of the input variables, and output a scalar through nonlinear transformation. Examples:

- sigmoid:

$$\sigma(v) = \frac{1}{1 + e^{-v}} = \frac{e^v}{e^v + 1}$$

- hyperbolic tangent (tanh):

$$\sigma(v) = \frac{e^v - e^{-v}}{e^v + e^{-v}}$$

- rectified linear unit (ReLU):

$$\sigma(v) = \max(0, v), \quad \text{soft approx.} \quad \ln(1 + e^v)$$

- And many others: exponential linear unit, arctangent, etc.

Overview

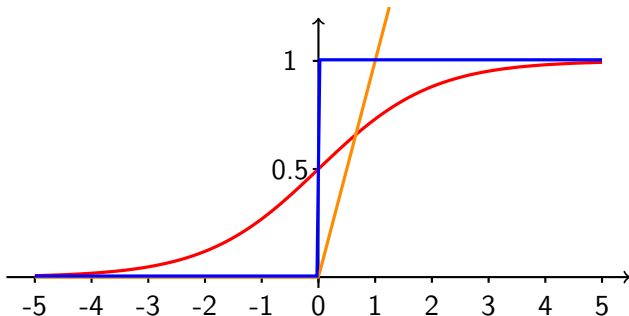
S.Lan

Deep Learning:  
Opportunities  
and Challenges

Course Overview

DNN Review

Computing



Sigmoid:  $(1 + e^{-v})^{-1}$

ReLU:  $\max(0, v)$ ,

Step function:  $I(v > 0)$

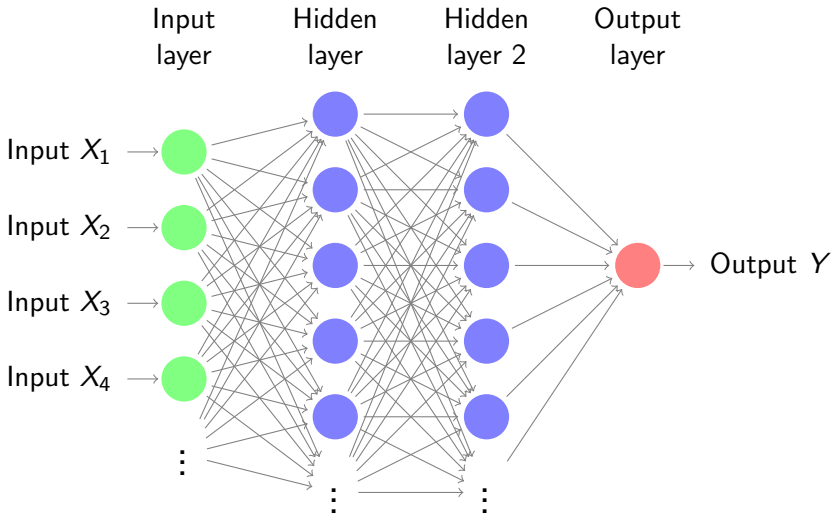
- Originally, a **step function**  $I(v > 0)$  was considered as the activation function (to mimic the biological interpretation). Hence for each neuron, signal is triggered only when  $x^T \alpha$  is above a certain threshold
- It was later recognized that the step function is **not smooth** enough for optimization, hence was replaced by a smoother threshold function, the sigmoid function
- “Feedforward” as signals can only pass to the next layer. There is no “cycle” in the model

## Universal Approximation Theorem (Cybenko, 1989; Hornik 1991)

Any continuous function  $f(x)$  on the space  $[0, 1]^p$  can be approximated (for any  $\epsilon > 0$ ) by a finite set of neurons with a bounded monotone-increasing activation function  $\sigma(\cdot)$ :

$$|f(x) - \sum_k w_k \sigma(\beta_k^T x + b_k)| < \epsilon$$

for some  $w_k$ ,  $\beta_k$ , and  $b_k$ . Hence, the functions defined by the neurons is dense.



- Try this a really cool website: <http://playground.tensorflow.org/>
- Implementation in [Python](#) :
  - packages: [sklearn.neural\\_network](#), ([TensorFlow](#), [PyTorch](#))
  - [MLPClassifier](#) implements a multi-layer perceptron (MLP) algorithm for classification.
  - [MLPRegressor](#) implements a multi-layer perceptron (MLP) for regression.
  - MLP trains using [Stochastic Gradient Descent](#) , [Adam](#) , or [L-BFGS](#) .
  - Important parameters:
    - number of neurons: [hidden\\_layer\\_sizes](#)
    - activation functions: [activation](#)
    - size of minibatches: [batch\\_size](#)
    - solver for back-propagation: [solver](#)
    - learning step sizes: [learning\\_rate](#), [learning\\_rate\\_init](#)
    - regularization: [alpha](#)



- The parameters (weights)  $\alpha$ 's and  $\beta$ 's need to be optimized.
- For a single hidden layer NN, we have

$$\{\alpha_1, \dots, \alpha_M\} : M(p+1) \text{ weights}$$

$$\{\beta_1, \dots, \beta_K\} : K(M+1) \text{ weights}$$

- where  $p$  is the number of non-intercept  $X$  features;  $M$  is the number of hidden neurons in a single layer; and  $K$  is the number of categories for classification.
- $K = 1$  if its a univariate regression problem.

- Neural Networks training is based on error minimization using a **Gradient Descent** algorithm, known as error **back-propagation**.
- For  $K$  classification, we minimize Deviance:

$$-\sum_{i=1}^n \sum_k^K \mathbf{1}\{y_i = k\} \log f_k(x_i)$$

- For univariate regression, we minimize RSS (since  $g$  is linear):

$$\sum_{i=1}^n (y_i - f(x_i))^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 \sigma(x^T \alpha_1) - \cdots \sigma(x^T \alpha_M))^2$$

- The objective function can be written as

$$R(\theta) = \sum_{i=1}^n R_i(\theta)$$

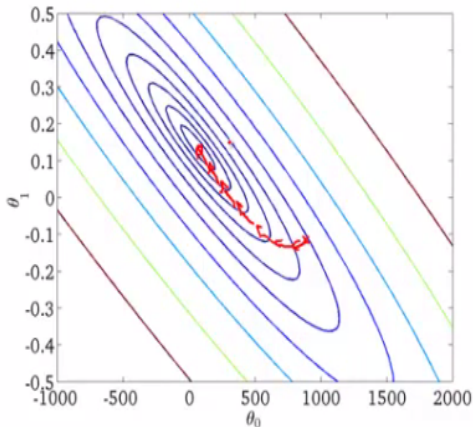
where  $R_i$  represents the deviance or residual sum of squares for the  $i$ th data point, and  $\theta$  represents an aggregated vector of all weights

- Initiate weights  $\theta^{(0)}$
- We then calculate the derivative wrt each of the weights evaluated at the current iteration value  $\theta^{(t)}$ :

$$\sum_{i=1}^n \frac{\partial R_i}{\beta_{km}} \bigg|_{\theta=\theta^{(t)}} \quad \sum_{i=1}^n \frac{\partial R_i}{\alpha_{mj}} \bigg|_{\theta=\theta^{(t)}}$$

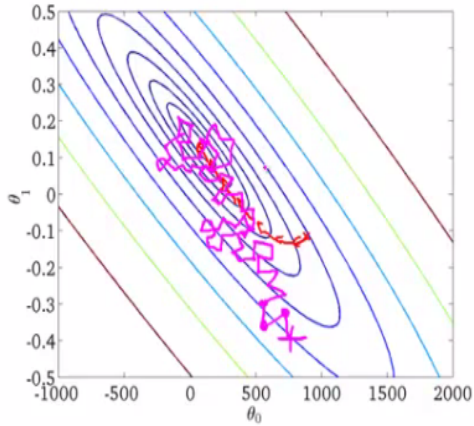
- **Stochastic GD**: the summation can be taken over a **random subset** of the  $n$  samples

Gradient Descent



vs.

Stochastic Gradient Descent



- The derivatives for  $K = 1$  regression case is essentially

$$\frac{\partial R_i}{\beta_m} = -2(y_i - f(x_i))z_{mi}$$

$$\frac{\partial R_i}{\alpha_{ml}} = -2(y_i - f(x_i))\beta_m \sigma'(\alpha_m^T x_i) x_{il}$$

- Some redundant calculations can be saved in the above equations. The property is called **back-propagation**.
- We then do the update, at the  $t$ -th iteration

$$\beta_m^{(t+1)} = \beta_m^{(t)} - \gamma \sum_{i=1}^n \frac{\partial R_i}{\beta_m^{(t)}}$$

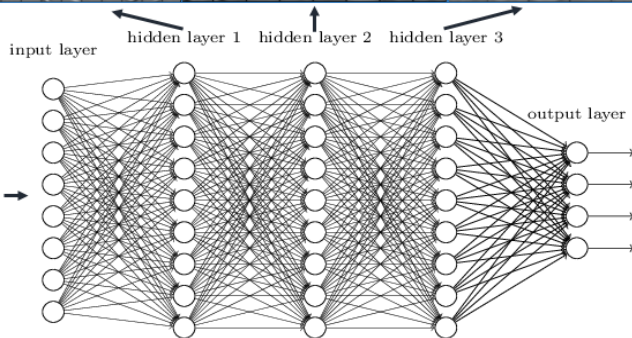
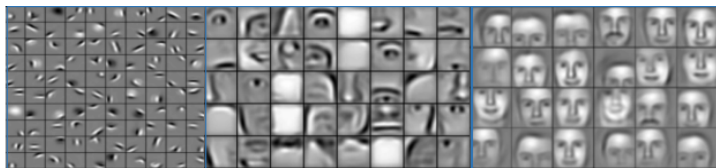
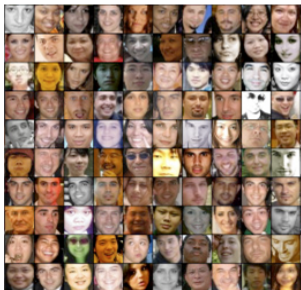
$$\alpha_{ml}^{(t+1)} = \alpha_{ml}^{(t)} - \gamma \sum_{i=1}^n \frac{\partial R_i}{\alpha_{ml}^{(t)}}$$

where  $\gamma$  is a step size for gradient descent.

- The derivatives can be calculated by Chain Rules
- The algorithm can be implemented by a forward-backward sweep over the network
- In the **forward** pass, compute the hidden variables and the output  $\hat{f}(x_i)$  based on the current weights  $\theta^{(t)}$
- In the **backward** pass, compute the derivatives, and update  $\theta^{(t)} \rightarrow \theta^{(t+1)}$

- **Deep Neural Networks** are one type of deep learning models.
- Deep neural Networks are just ... Neural Networks with more than one hidden layer.
- But neural networks have been around for more than 70 years... why it gets popular just in recent years?
  - computational issues
  - a better way to generate/construct features
  - ...

Deep neural networks learn hierarchical feature representations





Overview

S.Lan

Deep Learning:  
Opportunities  
and Challenges

Course Overview

DNN Review

Computing

1 Deep Learning: Opportunities and Challenges

2 Course Overview

3 DNN Review

4 Computing

- **Python** is a free software environment for statistical machine learning and graphics.  
— <https://www.python.org>
- **Conda** is an open source package / environment management system for Python.  
— <https://anaconda.org>
- **Jupyter** Notebook is a Python package to integrate code, equations, visualizations in document.  
— <https://jupyter.org>
- **Coding IDE** Both Eclipse (+PyDev) and Visual Studio Code are recommended.  
— <https://www.eclipse.org> (+ <https://www.pydev.org>) or <https://code.visualstudio.com>.
- **R Markdown** has also integrated Python. Check [link1](#) [link2](#).

- ASU Research Computing supports computing resources that can be used for fitting models.
  - <https://cores.research.asu.edu/research-computing>. One can start from [here](#).
- ASU has a dedicated website for Artificial Intelligence.
  - <https://ai.asu.edu>
- ASU provides multiple AI tools
  - <https://ai.asu.edu/ai-tools> You can get chatGPT-5.2 for free!
- Check out the AI resources and Guidelines
  - <https://instruction.thecollege.asu.edu/AIguidelines>