Google Classroom Code: mhxgl24

# Basics of Neural Networks

Deep Learning (DS-5006)
Dr. Adeel Mumtaz
Lecture 3
*Fall, 2022*
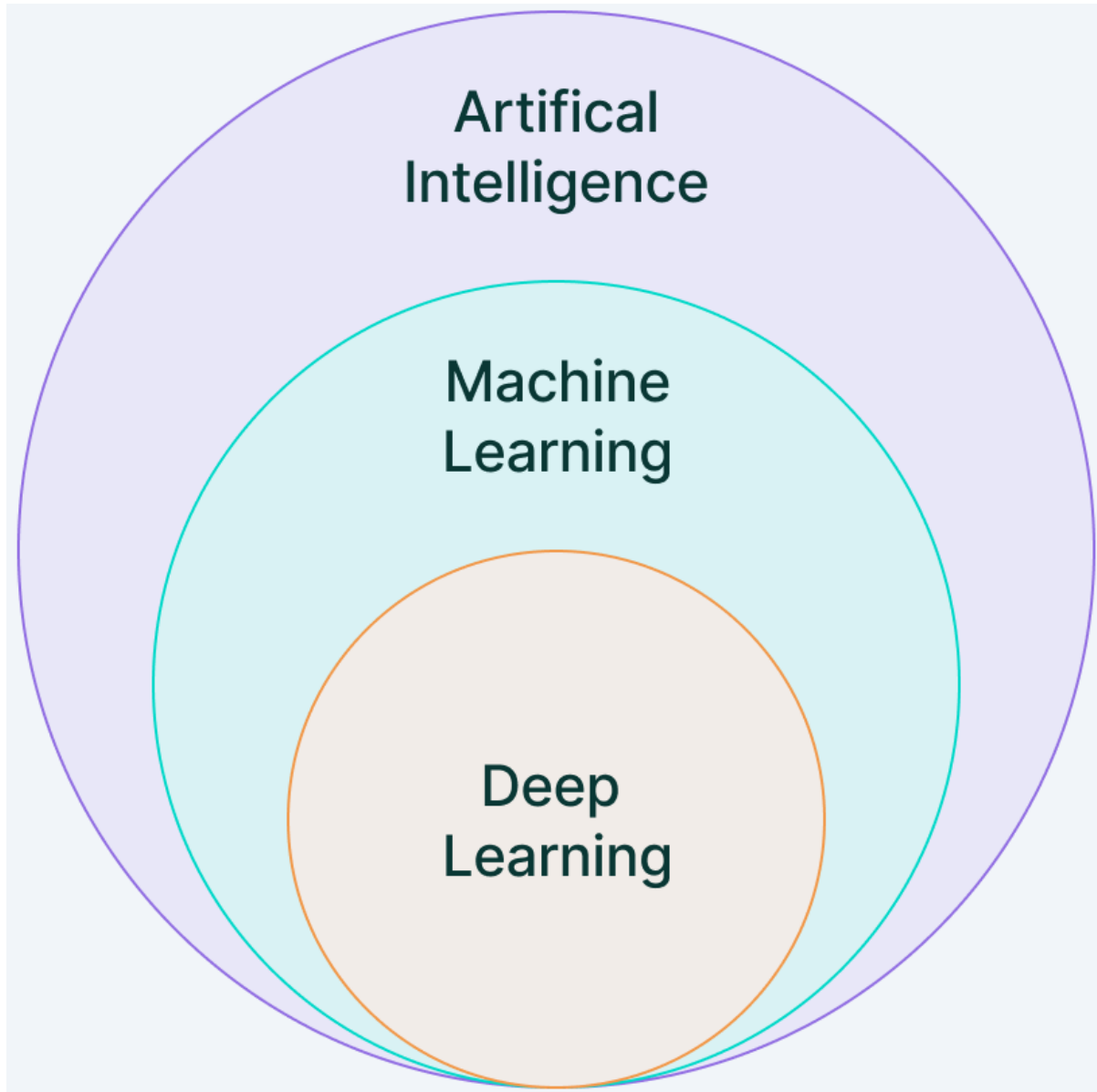
**National University**
**Of Computer and Emerging Sciences**
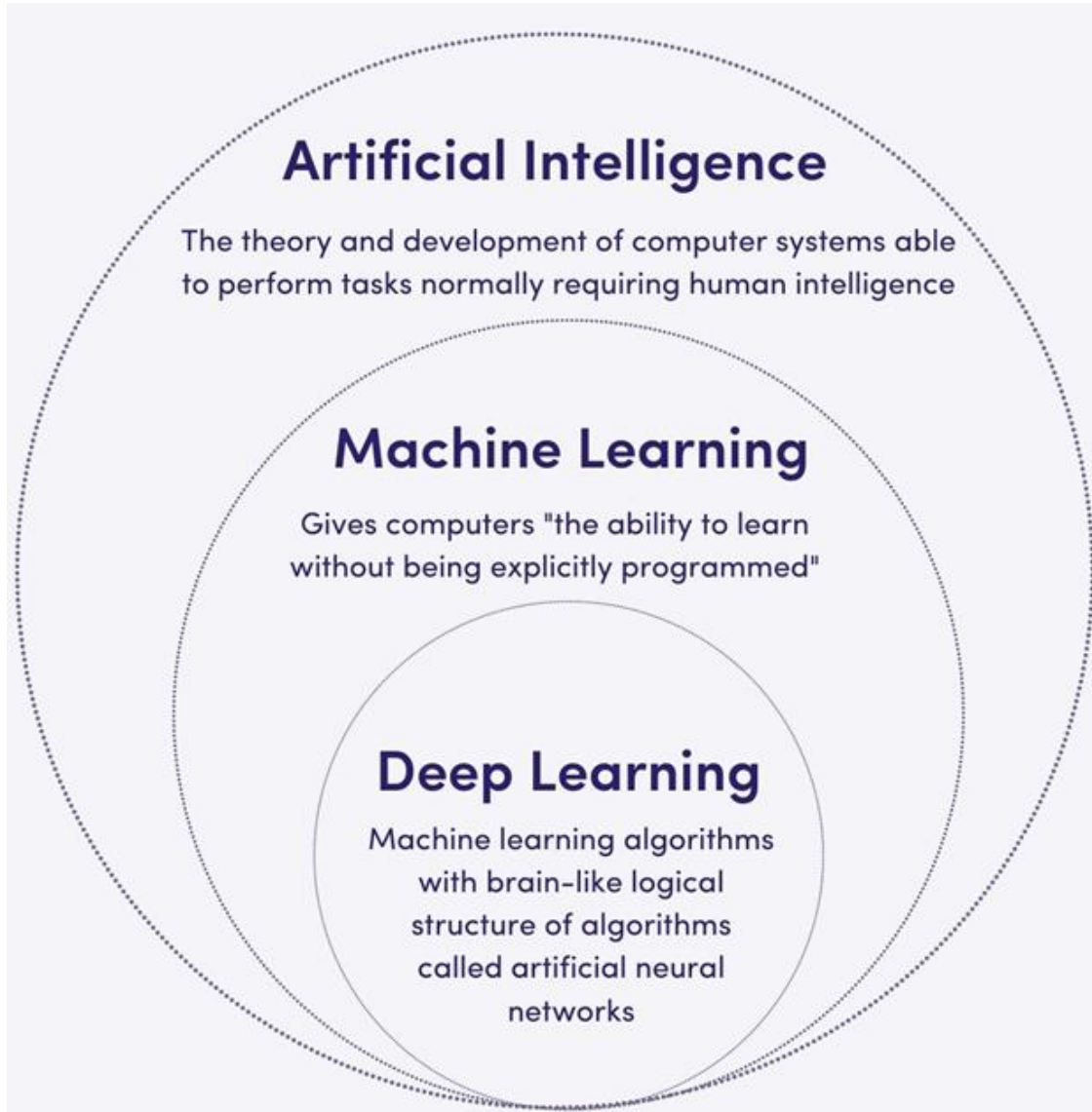
# Contents

- Basics of Neural Networks
  - Deep Learning
  - BNN Vs ANN
  - An Intuitive example (housing price prediction)
  - Structured and Unstructured data
  - Shallow Vs Deep Neural Networks
  - Deep Learning architectures
- A Simple Regression Problem (Theory)
  - Network Architecture
  - MSE Loss Function
  - Gradient Descent Algorithm
  - Learning Rate Effect
  - Training/Inference Loop
  - **Batch vs Stochastic vs Mini-Batch GD**

- A Simple Regression Problem (Numpy Implementation)
  - Data Generation
  - Data Splitting
  - Visualizing Data
  - Parameter initialization
  - Training Loop
    - Loss Calculation
    - Gradient Calculations
    - Parameter Updates
    - Validation Loss / Stopping Criteria
  - Plots/Training Curves
- NN Summary
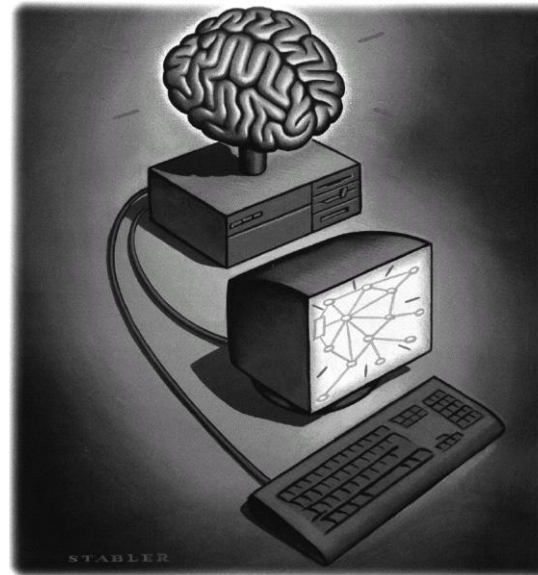- Home Task

# BASICS OF NEURAL NETWORKS

# Deep Learning

# Deep Learning

**Artificial Intelligence**

The theory and development of computer systems able to perform tasks normally requiring human intelligence

**Machine Learning**

Gives computers "the ability to learn without being explicitly programmed"

**Deep Learning**

Machine learning algorithms with brain-like logical structure of algorithms called artificial neural networks

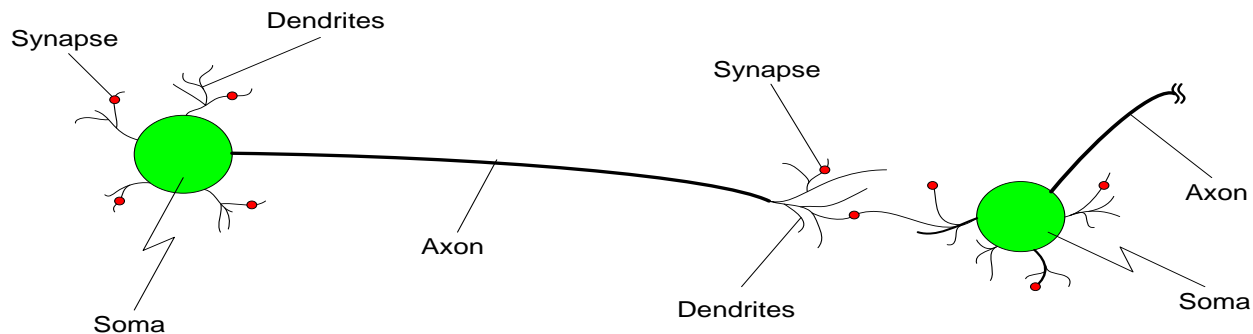- **The term Deep Learning refers to training very large Neural Network**

# What is a Neural Network

- Biologically motivated approach to machine learning

- Artificial neural network (ANN) is a machine learning approach that models human brain and consists of a number of artificial **neurons**

- Can be used for
  - Classification
  - Regression
  - Clustering
  - Association
  - Optimization
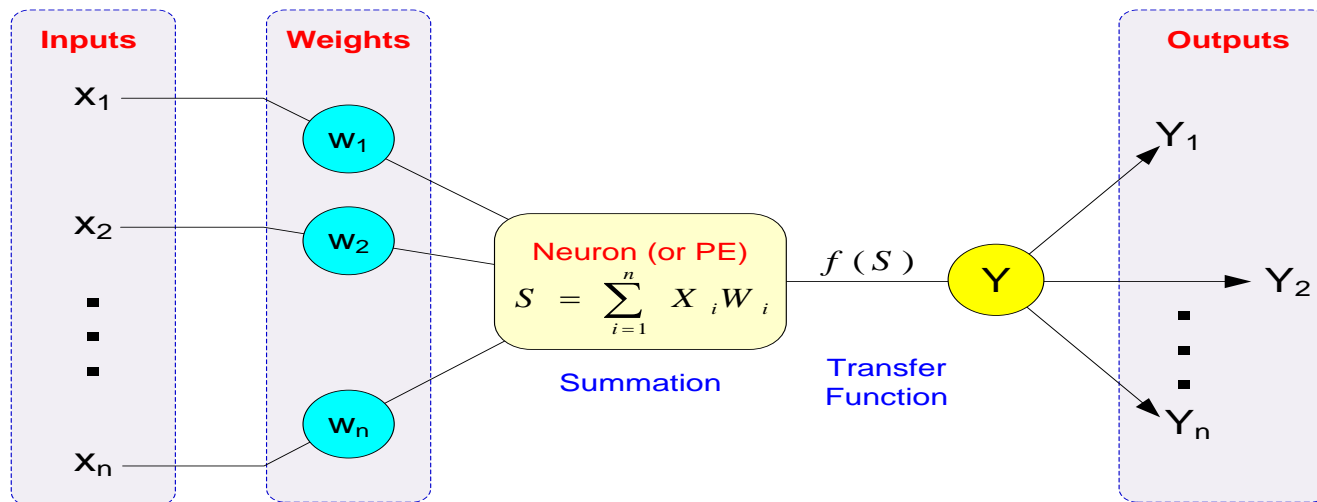


6

# Biological Neural Networks

- Fundamental processing element of a neural network is a ***neuron***
  - Dendrites, accept inputs
  - Soma, process the inputs
  - Axon, turns the processed inputs into outputs
  - Synapses, the electrochemical contacts between neurons
- A human brain has 100 billion neurons
- An ant brain has 250,000 neurons

**Two interconnected brain cells (neurons)**

# Artificial Neural Networks

- A single neuron (processing element – PE) with inputs and outputs
  - Inputs
  - Weights
  - Summation S or V
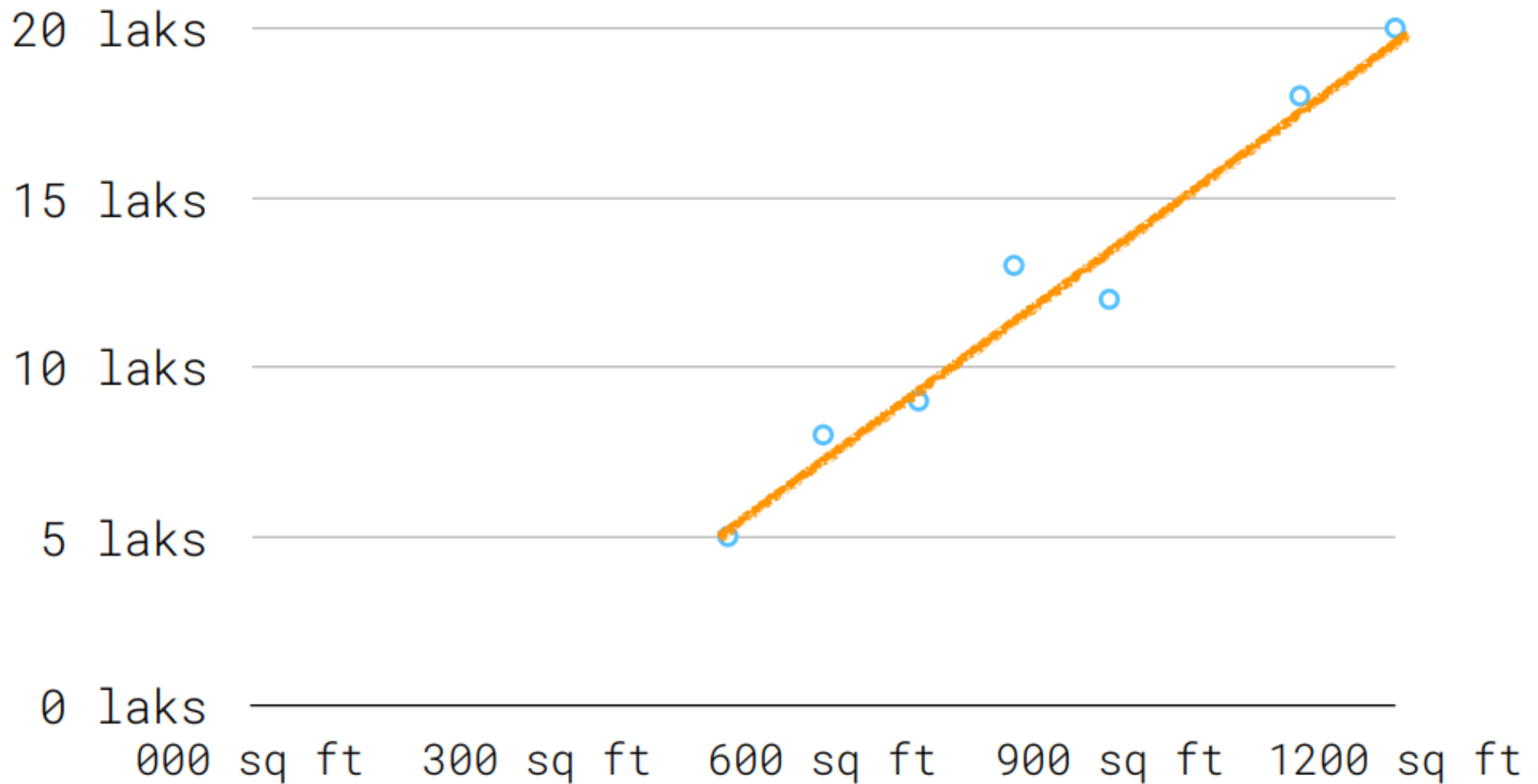  - Transfer function or Activation function, $\varphi \ or \ f$

**Inputs**

$x_1$
$x_2$
$\vdots$
$x_n$

**Weights**

$w_1$
$w_2$
$w_n$

Neuron (or PE)

$$S = \sum_{i=1}^{n} X_i W_i$$

Summation

$f(S)$

Y

Transfer Function

**Outputs**

$Y_1$
$Y_2$
$\vdots$
$Y_n$

# Comparison between ANN & BNN

| Biological | versus | Artificial | NNs |
|---|---|---|---|
| Soma | | Node | |
| Dendrites | | Input | |
| Axon | | Output | |
| Synapse | | Weight | |
| Slow | | Fast | |
| Many neurons ($10^9$) | | Few neurons ($\sim100$s) | |

# An Intuitive example (housing price prediction)

- Fit a function to predict the price of the house as a function of the size
  - linear regression?

| House Size (X) | Price (Y) |
| --- | --- |
| 500 | 5 Laks |
| 600 | 8 Laks |
| 700 | 9 Laks |
| 800 | 13 Laks |
| 900 | 12 Laks |
| 1100 | 18 Laks |
| 1200 | 20 Laks |

# Housing Price Prediction

# An Intuitive example (housing price prediction)

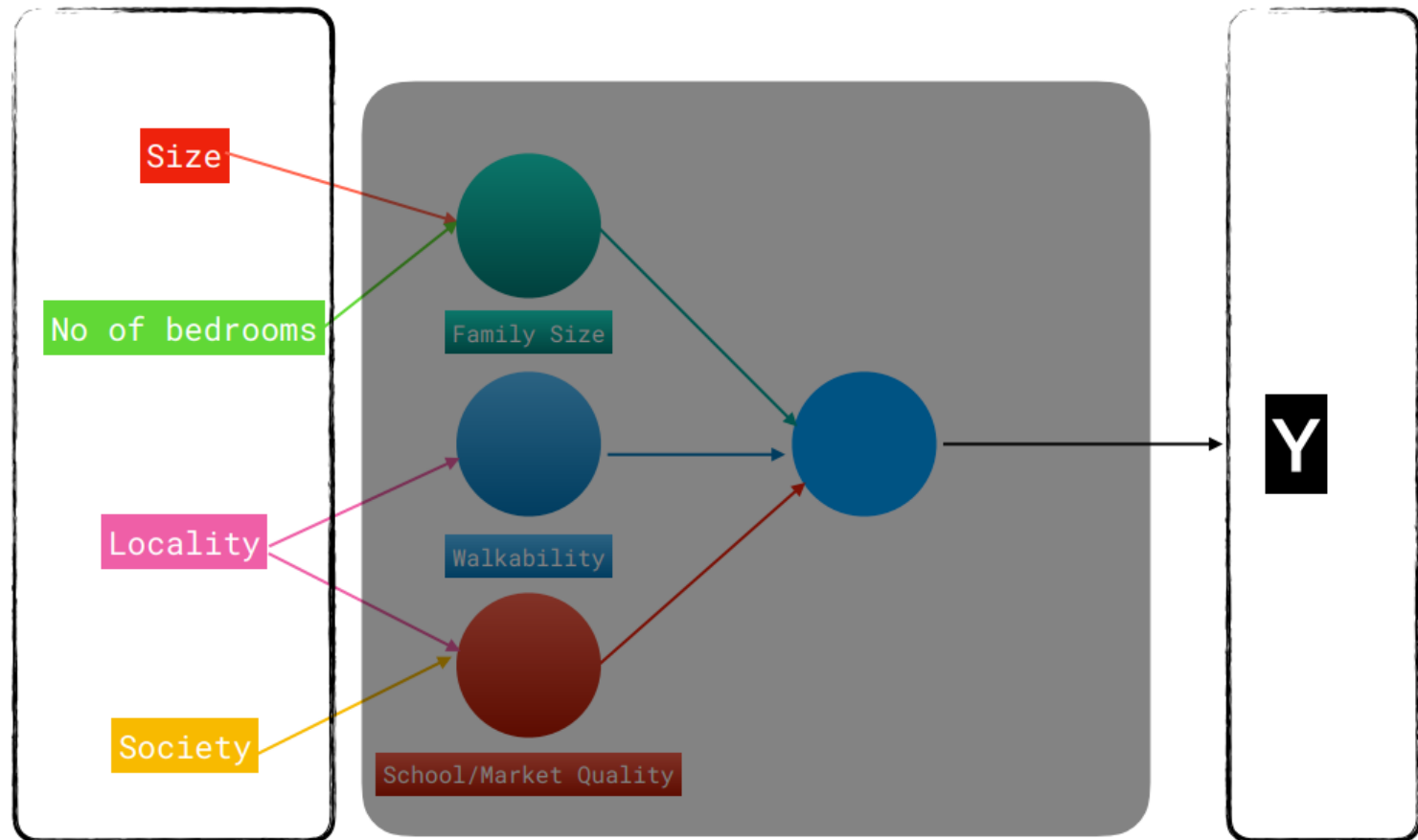- Price can't be negative

## ReLU Function

$$a = \max(0, z)$$

# An Intuitive example (housing price prediction)

- a very simple neural network
- Neuron takes inputs the size, computes the linear function, and then outputs the estimated price.
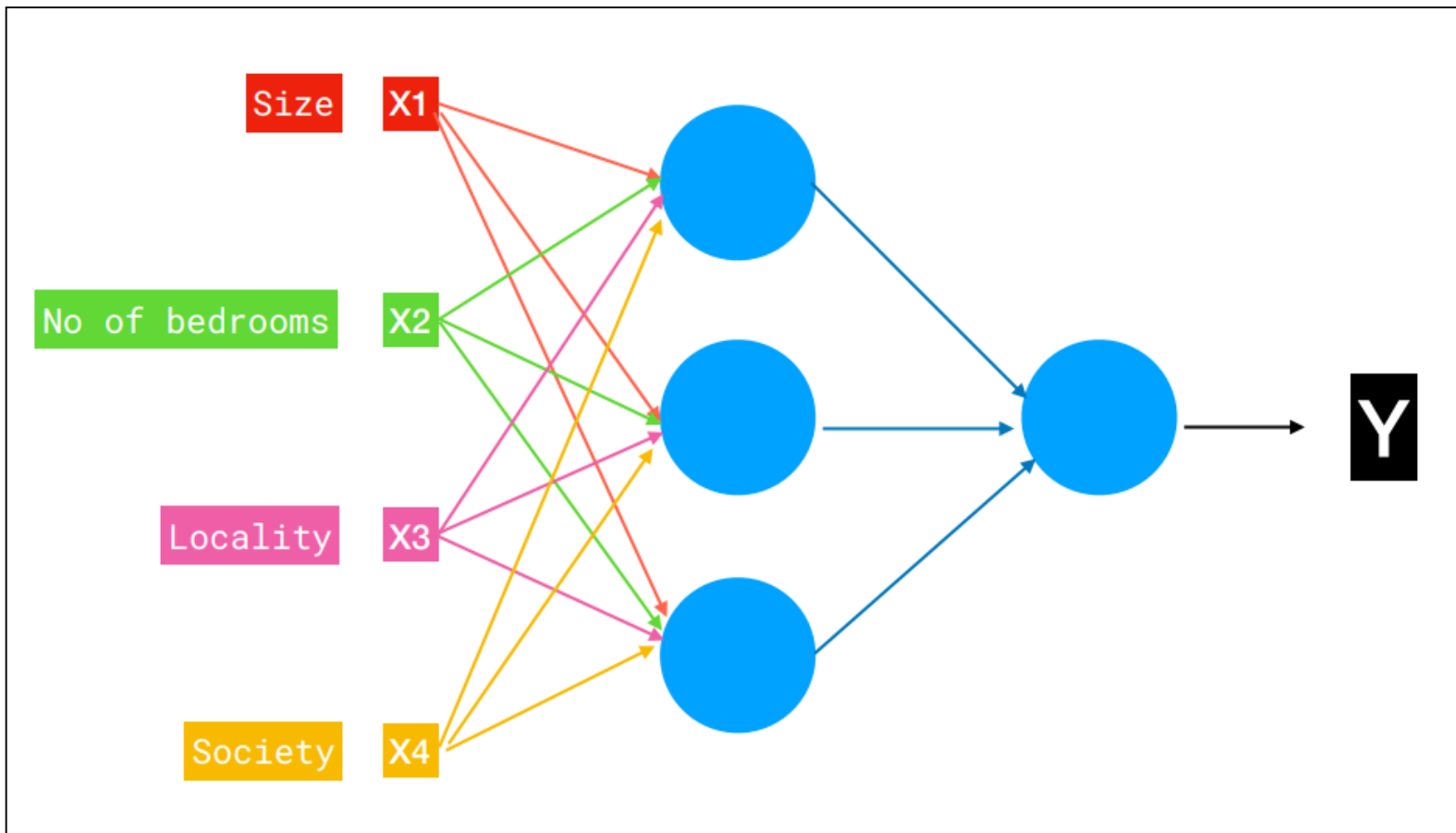
# An Intuitive example (housing price prediction)

- Other features that can affect the Price

# An Intuitive example (housing price prediction)

- A neural network with four inputs and one output
- One hidden layer with three hidden units
- Note that each of the hidden unit takes its inputs of all four features (Dense, Linear, Fully connected)

# Structured Data

| Size | #bedrooms | ... | Price (1000$s) |
|------|-----------|-----|----------------|
| 1200 | 2 | | 300 |
| 1500 | 3 | | 400 |
| 2000 | 3 | | 480 |
| ⋮ | ⋮ | | ⋮ |
| 3000 | 4 | | 520 |

| User Region | Ad Id | ... | Ad revenue($) |
|-------------|-------|-----|---------------|
| USA | 1005 | | 0.5 |
| UK | 1009 | | 0.3 |
| USA | 998 | | 0.5 |
| ⋮ | ⋮ | | ⋮ |
| CAN | 2104 | | 0.45 |

# Unstructured Data



Audio

Image

Once upon a time in the history of…

Text

# Unstructured Data

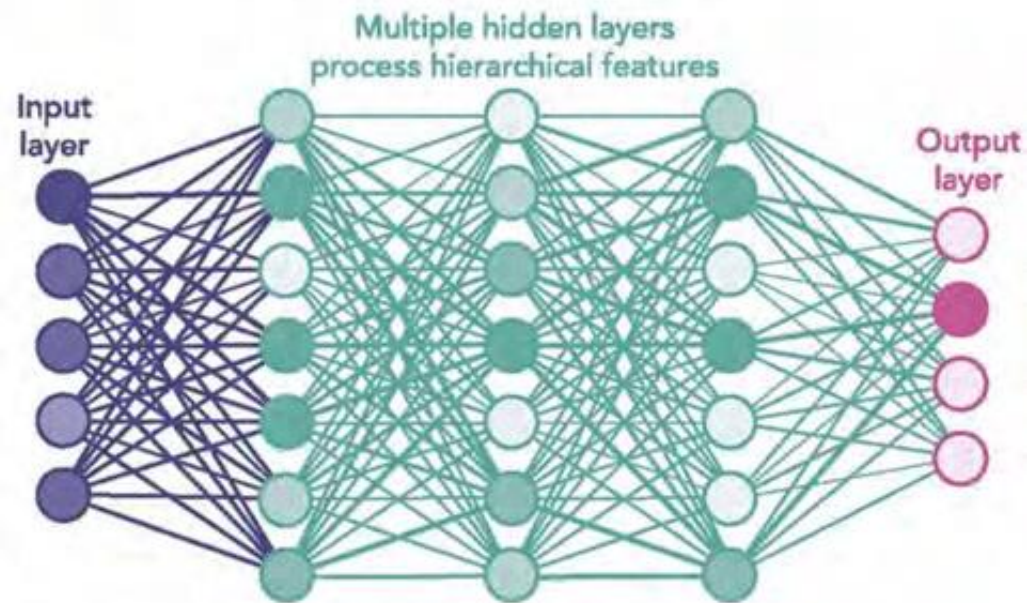| | | | |
|---|---|---|---|
| **Text files and documents** | **Server, website and application logs** | **Sensor data** | **Images** |
| **Video files** | **Audio files** | **Emails** | **Social media data** |

**Deep Neural Networks are very good at processing these Unstructured data**

# Shallow Vs Deep Neural Networks

# Shallow Vs Deep Neural Networks

| Factors | Shallow Neural Network (SNN) | Deep Neural Network (DNN) |
|---|---|---|
| Feature Engineering | 1. Individual feature extraction process is required. Various features cited in the literature are histogram oriented gradients, speeded up robust features, and local binary patterns. | 1. Replace the hand-crafted features and directly work on the entire input. Thus, more practical for complex datasets. |
| Data Size Dependency | 2. Needs a lesser quantity of data. | 2. Needs vast volumes of data. |

# Deep Learning Architectures

- Universal Function Approximators
- The number of trainable parameters increases drastically with an increase in the size of the image



**Artificial Neural Network (ANN or MLP)**

# Deep Learning Architectures

- CNN learns the filters automatically
- CNN captures the spatial features from an image, identify object location accurately
- CNN follows the concept of parameter sharing



**Convolutional Neural Network (CNN)**

# Deep Learning Architectures

- RNN captures the sequential information present in the input data
- RNNs share the parameters across different time steps



**Recurrent Neural Network (RNN)**

# Deep Learning Architectures

| | MLP | RNN | CNN |
|---|---|---|---|
| Data | Tabular data | Sequence data (Time Series, Text, Audio) | Image data |
| Recurrent connections | No | Yes | No |
| Parameter sharing | No | Yes | Yes |
| Spatial relationship | No | No | Yes |
| Vanishing & Exploding Gradient | Yes | Yes | Yes |

# Deep Learning Architectures



**Attention Based Neural Network (Transformer)**

# Deep Learning Architectures



**Generative Adversarial Network (GAN)**

# A SIMPLE REGRESSION PROBLEM (THEORY)

# Simple Linear Regression

$$\hat{y} = b + wX$$

Dependent Variable/ output

y-intercept/ bias constant

Coefficient of regression

Features or Independent variable

# Simple Linear Regression



regression line
$$y = b + wX$$

data points

# Simple Linear Regression

# Simple Linear Regression

$$Cost(b,w) = \mathcal{L}(b,w) = \frac{1}{N}\sum_{i=1}^{N}(\widehat{y}_i - y_i)^2$$

$$= \mathcal{L}(x,b,w) = \frac{1}{N}\sum_{i=1}^{N}(wx_i + b - y_i)^2$$

$Cost(b,w)$

# Simple Linear Regression



The Linear Unit: $y = wx + b$

- **NN=Architecture + Parameters**
- **Training NN = Given data learn best Parameters which gives minimum loss (usually an iterative process/algorithm)**

# Simple Linear Regression

# Gradient Descent Algorithm

$$W := W - \alpha \frac{\partial}{\partial W} cost(W)$$

# Simple Linear Regression

$$\mathcal{L}(x, b, w) = \frac{1}{N} \sum_{i=1}^{N} (wx_i + b - y_i)^2$$

$$dw = \frac{\partial \mathcal{L}(x, b, w)}{\partial w} = 2 * \frac{1}{N} \sum_{i=1}^{N} (wx_i + b - y_i)(x_i)$$

$$= 2 * \frac{1}{N} \sum_{i=1}^{N} (error_i)(x_i) = 2 * mean(error * x)$$

$$db = \frac{\partial \mathcal{L}(x, b, w)}{\partial b} = 2 * \frac{1}{N} \sum_{i=1}^{N} (error_i) = 2 * mean(error)$$

$$w_{new} = w_{old} - \alpha * dw$$

$$b_{new} = b_{old} - \alpha * db$$

# Effect of Learning Rate



**Too low**

$J(\theta)$

$\theta$

A small learning rate requires many updates before reaching the minimum point

**Just right**

$J(\theta)$

$\theta$

The optimal learning rate swiftly reaches the minimum point

**Too high**

$J(\theta)$

$\theta$

Too large of a learning rate causes drastic updates which lead to divergent behaviors

# Effect of Learning Rate



loss

very high learning rate

low learning rate

high learning rate

good learning rate

epoch

# Training Loop

**Save Model**
$w \ \& \ b$

YES

**Load Training & Validation Data**
$(x_i, y_i)$

NO

**Stopping Criteria Validation Loss?**

**Randomly Initialize Parameters**
$(w, b)$
**Choose a learning rate:** $\alpha$

**Optimizer**

**Parameter Update**

$w = w_{old} - \alpha * dw$
$b = b_{old} - \alpha * db$

**Forward Pass (using current** $w, \boldsymbol{b}$)
$\widehat{y}_i = wx_i + b$

**Compute/Plot  Loss Value**

$$\mathcal{L}(x, b, w) = \frac{1}{N} \sum_{i=1}^{N} (\widehat{y}_i - y_i)^2$$

**Backward Pass (Gradients)**

$dw = 2 * mean(error * x)$
$db = 2 * mean(error)$ where,
$error_i = \widehat{y}_i - y_i$

38

# Inference/Testing

Load Test Data
$(x_i, y_i)$

Load Saved Model
$(w, b)$

Forward Pass (using **saved** $\boldsymbol{w}, \boldsymbol{b}$)
$\widehat{y}_i = wx_i + b$

Compute/Plot Loss Value

$$\mathcal{L}(x, b, w) = \frac{1}{N} \sum_{i=1}^{N} (\widehat{y}_i - y_i)^2$$

# Batch Gradient Descent

- An epoch refers to one cycle through the full training dataset
- All the training data is taken into consideration to update parameters
- Uses mean of gradients to update parameters
- One step/iteration of gradient descent in one epoch
- Great for convex or relatively smooth error surfaces
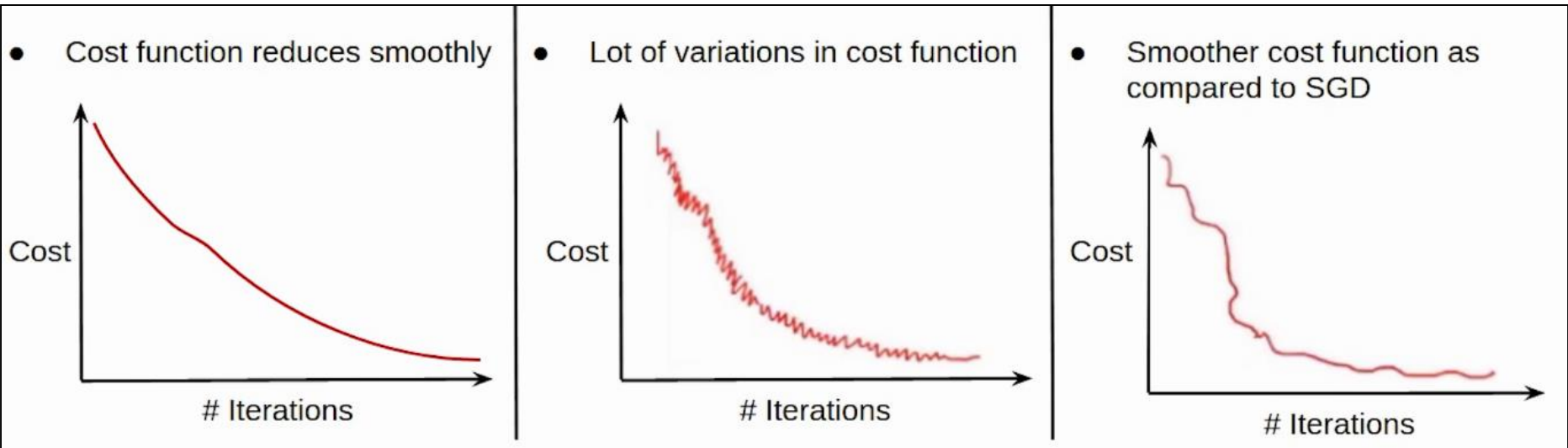- The graph of cost vs epochs is also quite smooth

# Stochastic Gradient Descent

- If our dataset has 5 million examples, then just to take one step the model will calculate the gradients 5 million examples (<span style="color:red">inefficient</span>)

- In Stochastic Gradient Descent (SGD), we consider just <span style="color:blue">one example</span> at a time to take a single step.

- Cost will fluctuate over the training examples

- N <span style="color:blue">steps/iterations</span> of gradient descent in one epoch

# Mini Batch Gradient Descent

- We use a batch of a fixed number of training examples which is less than the actual dataset and call it a mini-batch

- Average cost over the epochs in mini-batch gradient descent fluctuates because we are averaging a small number of examples at a time.

- N/batch_size steps/iterations of gradient descent in one epoch

# Batch vs Stochastic vs Mini-Batch



| Batch GD | Stochastic GD | Mini-Batch GD |
| --- | --- | --- |
| • Cost function reduces smoothly | • Lot of variations in cost function | • Smoother cost function as compared to SGD |

**We can divide the dataset of 2000 examples into batches of 500 then it will take ? iterations to complete 1 epoch.**

# A SIMPLE REGRESSION PROBLEM (NUMPY IMPLEMENTATION)

# Dataset Generation

```python
import numpy as np
import matplotlib.pyplot as plt

#data generation
true_w=2
true_b=1
N=100

np.random.seed(100)
#get N uniformly distributed values
x=np.random.rand(N,1)
#get N noise values from standard normal distribution
epsilon=0.1*np.random.randn(N,1)
y=true_w*x+true_b+epsilon
```

# Splitting Dataset into train/validation sets

```python
#Splitting Data into train and validation
idx=np.arange(N)
np.random.shuffle(idx)
idx_train=idx[:int(0.8*N)]
idx_test=idx[int(0.8*N):]
x_train, y_train = x[idx_train],y[idx_train]
x_val, y_val = x[idx_test],y[idx_test]
```
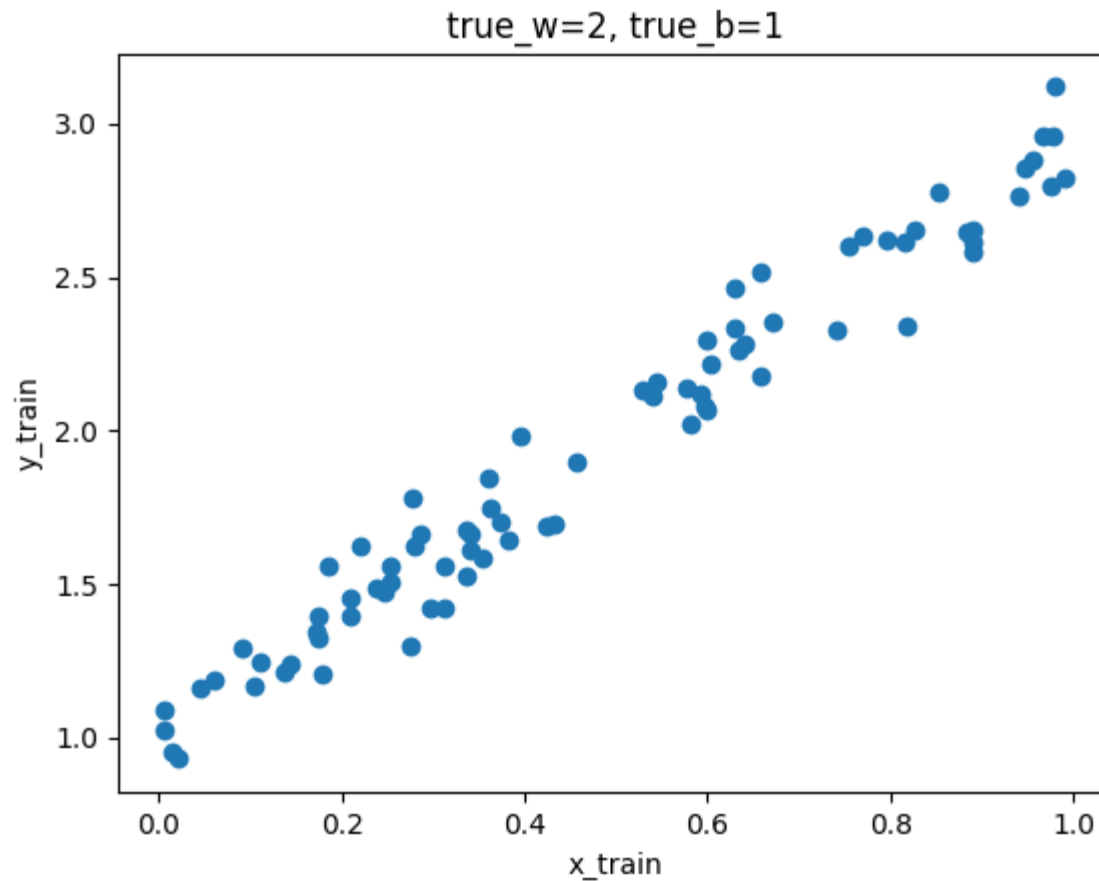
```
∨ WATCH
  > x_val.shape: (20, 1)
  > y_val.shape: (20, 1)
  > x_train.shape: (80, 1)
  > x_train.shape: (80, 1)
```
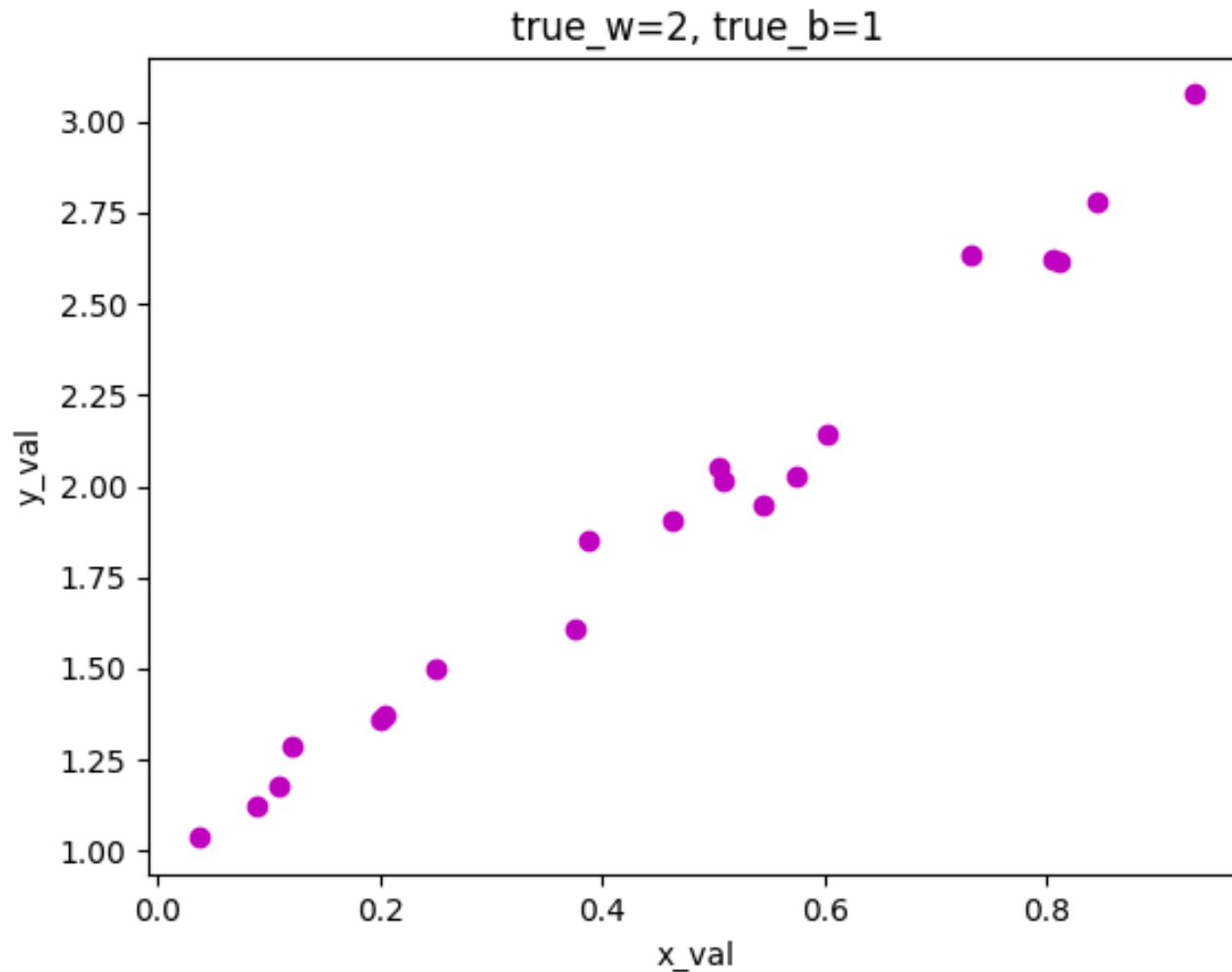
# Splitting Dataset into train/validation sets

```
25    #plotting tain and val data
26    plt.figure('1')
27    plt.scatter(x_train,y_train)
28    plt.xlabel('x_train')
29    plt.ylabel('y_train')
30    plt.title(f'true_w={true_w}, true_b={true_b}')
31    plt.figure('2')
32    plt.scatter(x_val,y_val,color = 'm')
33    plt.xlabel('x_val')
34    plt.ylabel('y_val')
35    plt.title(f'true_w={true_w}, true_b={true_b}')
36    plt.show(block=True)
```

# Visualizing Datasets



true_w=2, true_b=1

# Visualizing Datasets



true_w=2, true_b=1

# Training Loop

```python
#training loop
#initializing parameters
trainLosses=[]
valLosses=[]
lr=0.1
w=np.random.randn(1)
b=np.random.randn(1)
for i in range(100):
    #forward pass
    yhat=w*x_train+b #note vectorized operation
    #MSE loss
    error=yhat-y_train
    loss= (error**2).mean()
    trainLosses.append(loss)
    #computing gradients
    db=2*error.mean()
    dw=2*(x_train*error).mean()
    #weight update
    b=b-lr*db
    w=w-lr*dw
```

# Validation Loss

```python
#val MSE loss
yhatVal=w*x_val+b
errorVal=yhatVal-y_val
valLoss= (errorVal**2).mean()
valLosses.append(valLoss)

#stopping condition
if(valLoss<0.0001):
    break

print(f'train loss={loss}, val loss={valLoss}, w={w}, b={b}')
```
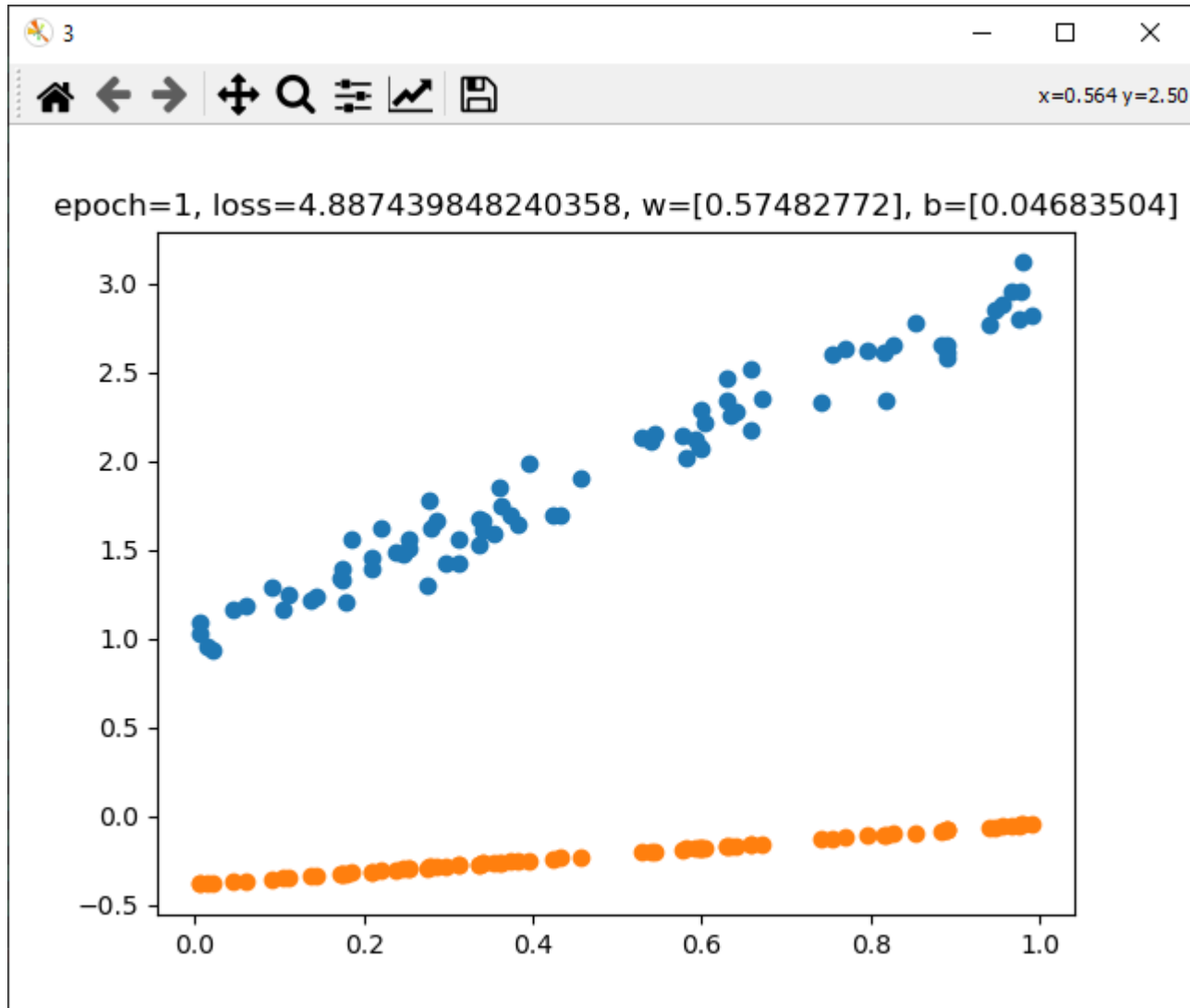
# Plots (Regression Fitting)

```python
72    #training data plot
73    plt.figure('3')
74    plt.cla()
75    plt.scatter(x_train,y_train)
76    plt.scatter(x_train,yhat)
77    plt.title(f'epoch={i}, loss={loss}, w={w}, b={b}')
78    plt.show(block=False)
79    plt.pause(1)
80
81    #validation data plot
82    plt.figure('4')
83    plt.cla()
84    plt.scatter(x_val,y_val)
85    plt.scatter(x_val,yhatVal)
86    plt.title(f'epoch={i}, ValLoss={valLoss}, w={w}, b={b}')
87    plt.show(block=False)
88    plt.pause(1)
```
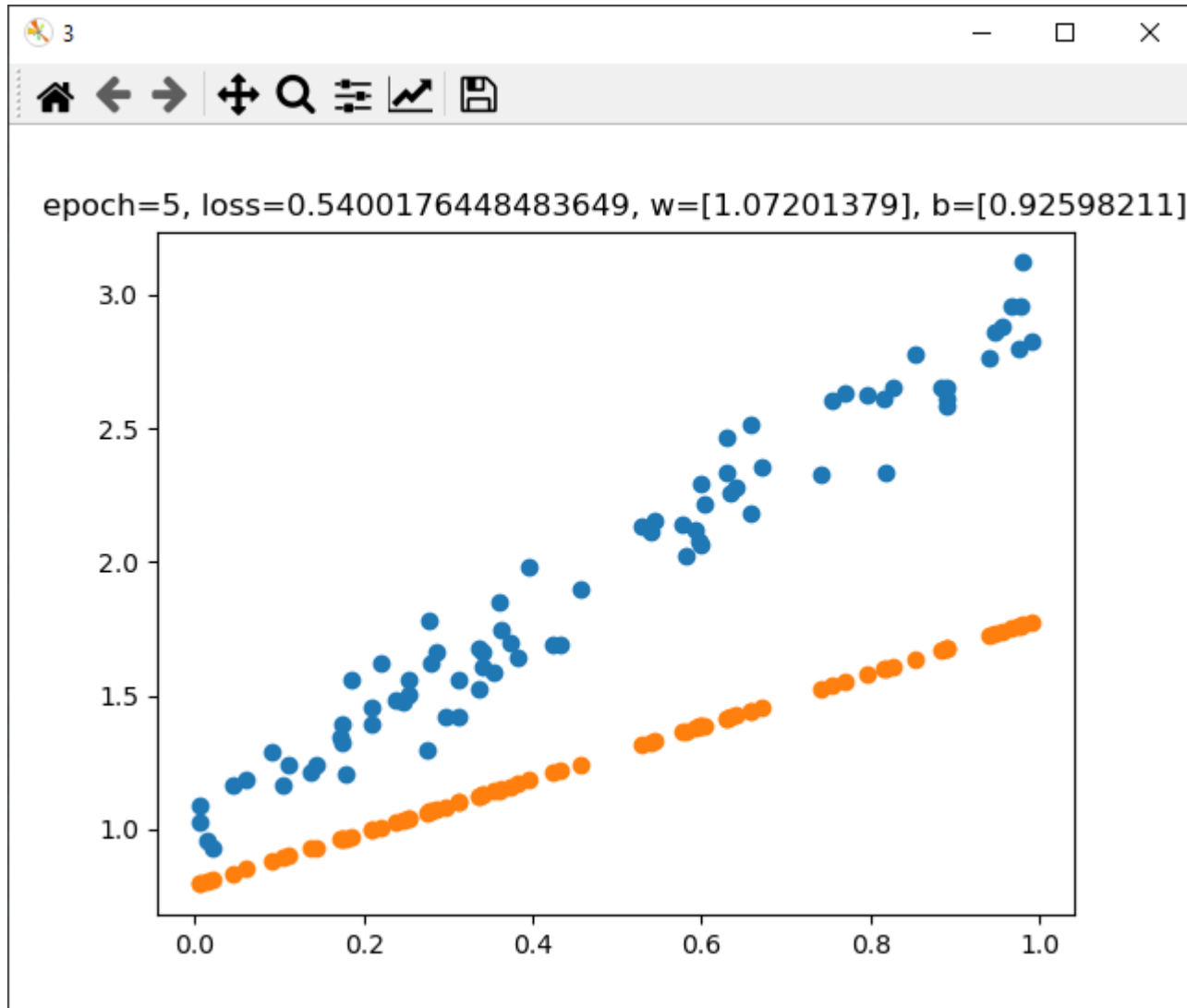
# Plots (Loss vs Epoch)

```
 90        #trainLoss vs Epoch
 91        plt.figure('5')
 92        plt.cla()
 93        plt.plot(trainLosses)
 94        plt.xlabel('Epoch')
 95        plt.ylabel('trainLoss')
 96        plt.title(f'Training Loss Vs Epoch')
 97        plt.show(block=False)
 98        plt.pause(1)
 99
100        #validationLoss vs Epoch
101        plt.figure('6')
102        plt.cla()
103        plt.plot(valLosses,color='m')
104        plt.xlabel('Epoch')
105        plt.ylabel('valLoss')
106        plt.title(f'Validation Loss Vs Epoch')
107        plt.show(block=False)
108        plt.pause(1)
```
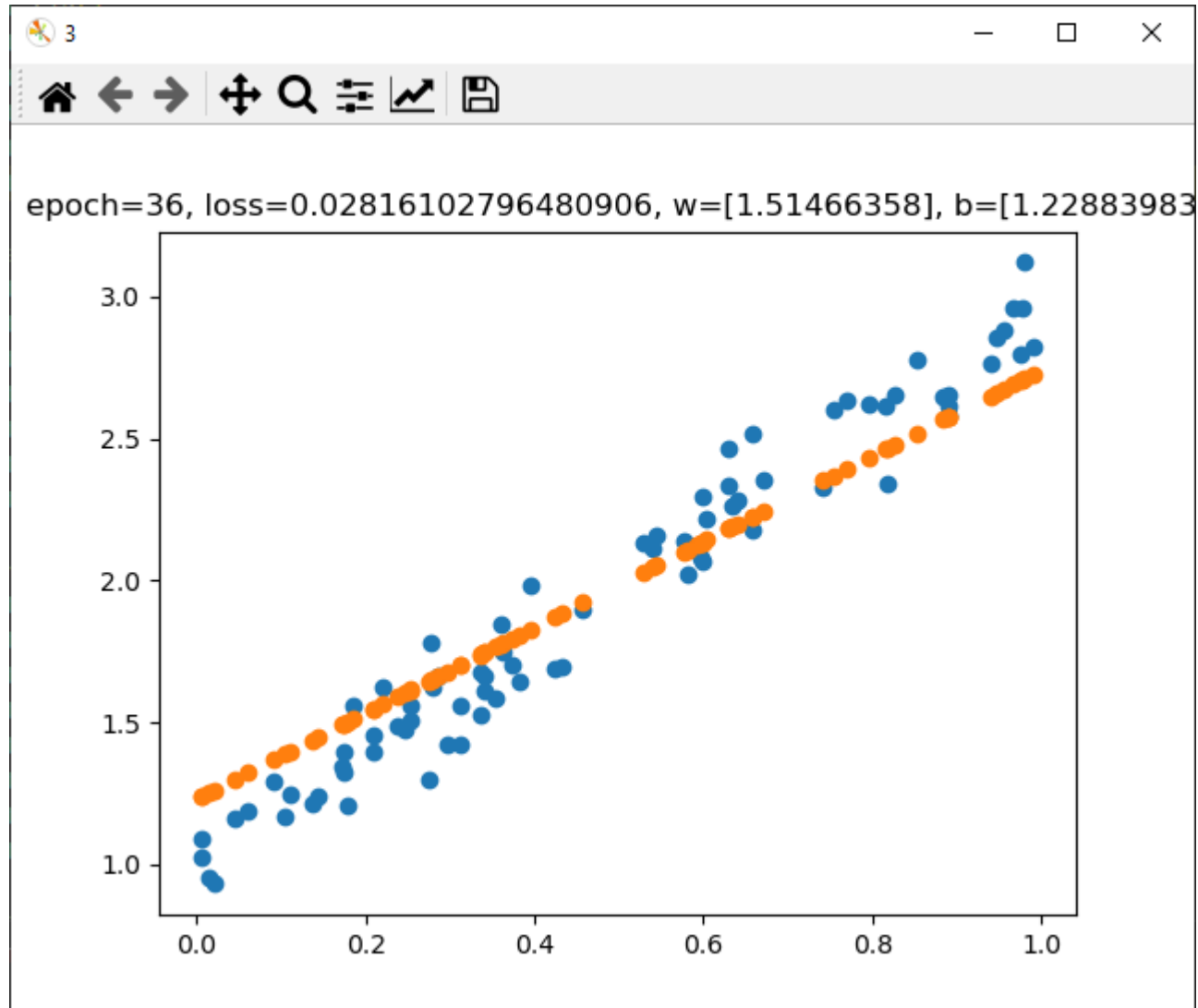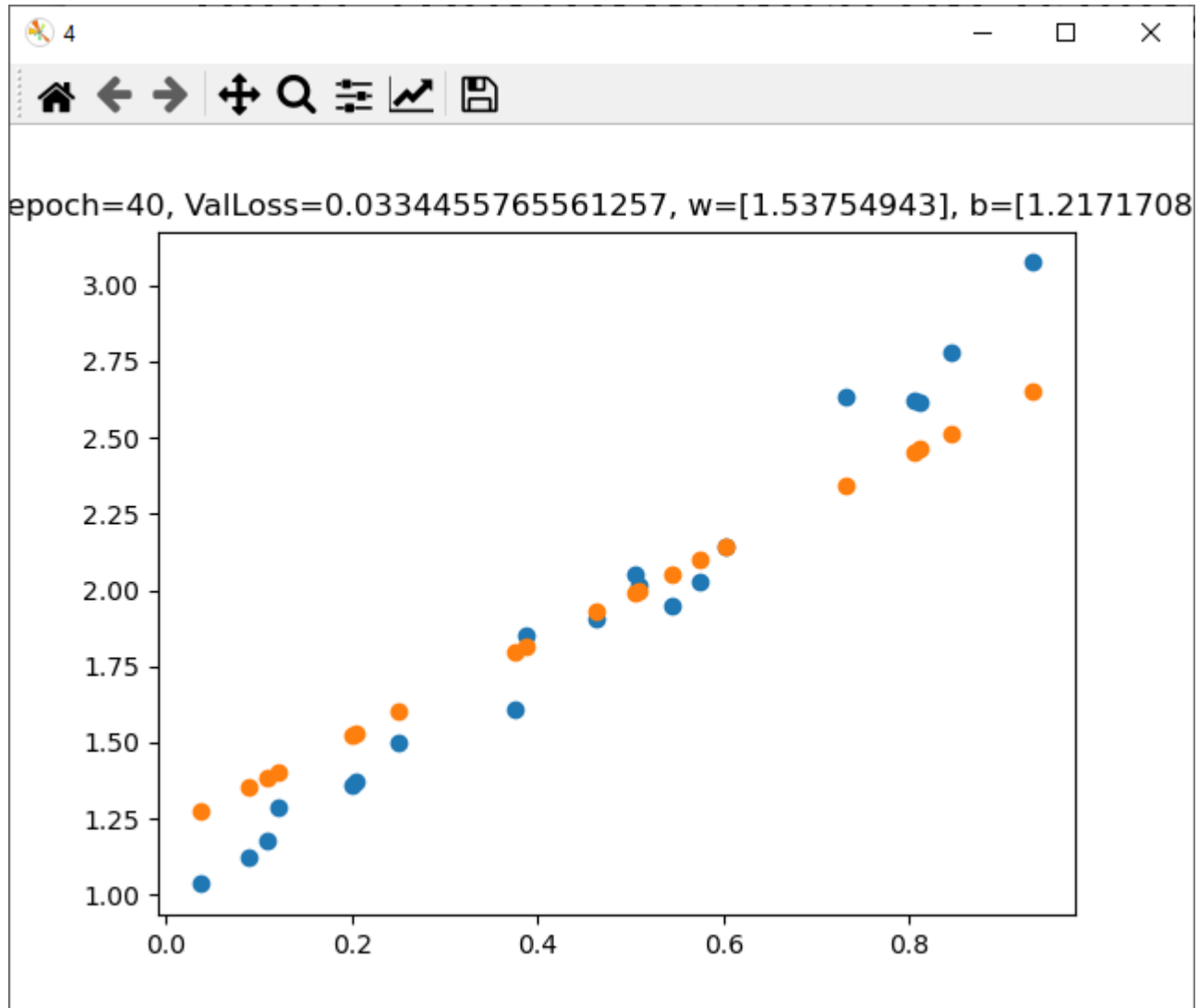
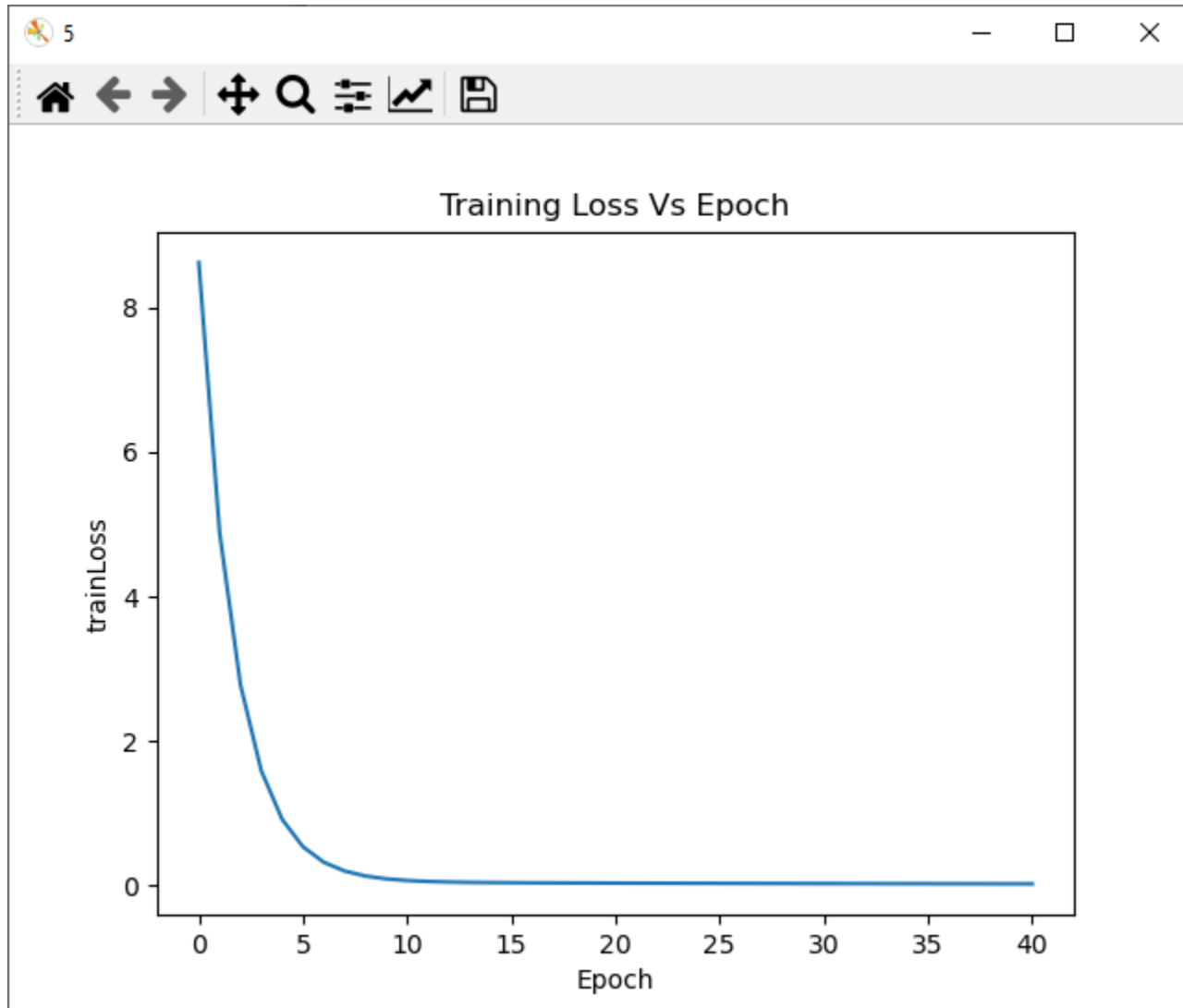# Plots (Convergence of regression)

# Plots (Convergence of regression)



epoch=5, loss=0.5400176448483649, w=[1.07201379], b=[0.92598211]

# Plots (Convergence of regression)



epoch=36, loss=0.02816102796480906, w=[1.51466358], b=[1.22883983

# Plots (Performance on Validation Set)



epoch=40, ValLoss=0.0334455765561257, w=[1.53754943], b=[1.2171708

# Epoch vs Train Loss

# Epoch vs Validation Loss

# NN Summary

- Data Set, Training, Validation, Test
- Cost/Loss Function
  - MSE Loss for regression
- Training Loop
- Optimizer
- Parameters
- Learning Rate
- Epoch
- Batch
- Loading/Saving Model

# Home Task

- Compare learning curves for different values of learning rate

- Convert code from Batch Gradient Descent to Stochastic Gradient Descent and compare learning curves