

# Neural networks for regression

Abhinav Pradeep

October 22, 2024

# Function approximation

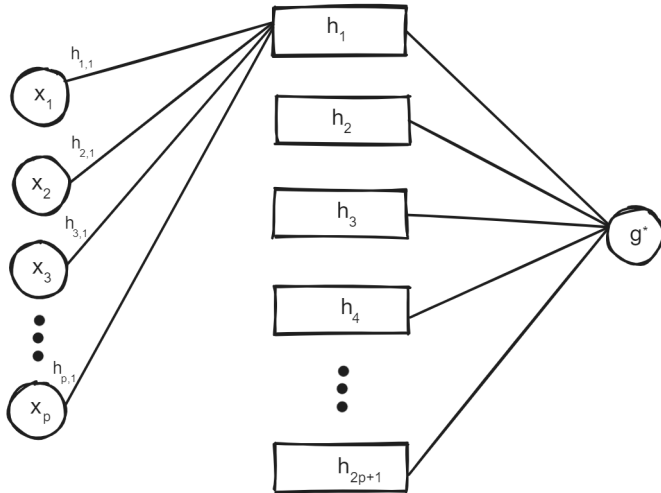
Theorem 9.1 from the textbook:

For  $p \geq 2$  every continuous function  $g^*(\mathbf{x}) : [0, 1]^p \mapsto \mathbb{R}$  can be written as:

$$g^*(\mathbf{x}) = \sum_{j=1}^{2p+1} h_j \left( \sum_{i=1}^p h_{ij}(x_i) \right)$$

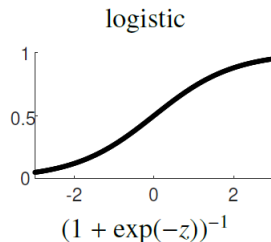
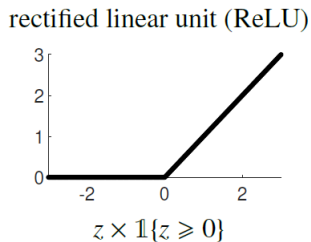
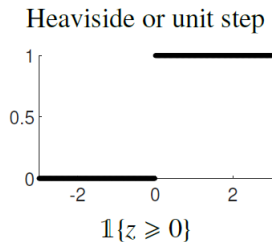
Where  $\{h_j, h_{ij}\}$  are continuous univariate functions. So every continuous multivariate function can be written as a sum of  $p(2p + 1)$  univariate functions.

# Function approximation

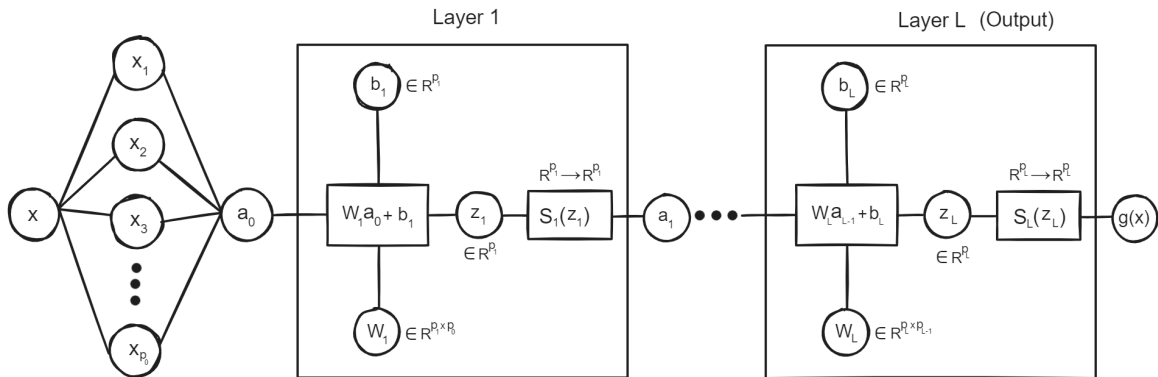


# Function approximation

We do not know the set of  $p(2p + 1)$  functions  $\{h_j, h_{ij}\}$  and they are likely non-linear. We can replace the by a larger set of functions. Call these activation functions  $S(z)$ . Examples are (from the textbook):



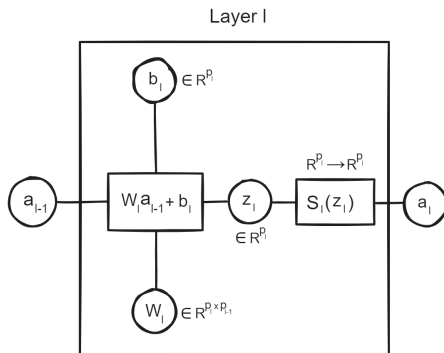
# Feed-forward networks



# Feed-forward networks

The activation function  $\mathbf{a}_l = \mathbf{S}_l(\mathbf{z}_l)$  for the hidden layers ( $l = 1, 2 \dots L - 1$ ) is

$\mathbf{S}_l(\mathbf{z}) = \begin{bmatrix} S(z_1) & S(z_2) & \dots & S(z_{p_l}) \end{bmatrix}$ . The output layers activation  $\mathbf{S}_L$  function is more general and depends on the goal.



# Feed-forward propagation

---

## Algorithm 1 Feed-forward propagation

---

**Input:**  $\mathbf{x} \in \mathbb{R}^{p_0}$ ,  $\{\mathbf{W}_l \in \mathbb{R}^{p_l \times p_{l-1}}\}$ ,  $\{\mathbf{b}_l \in \mathbb{R}^{p_l}\}$ ,  $\{\mathbf{S}_l\}$

**Output:**  $g(\mathbf{x})$

- 1:  $\mathbf{a}_0 \leftarrow \mathbf{x}$
  - 2: **for**  $l = 1$  to  $L$  **do**
  - 3:      $\mathbf{z}_l \leftarrow \mathbf{W}_l \mathbf{a}_{l-1} + \mathbf{b}_l$
  - 4:      $\mathbf{a}_l \leftarrow \mathbf{S}_l(\mathbf{z}_l)$
  - 5: **end for**
  - 6: **return**  $g(\mathbf{x}) \leftarrow \mathbf{a}_L$
-

# Back propagation

For steepest descent type algorithms we wish to compute the gradient of training loss. First we gather all the parameters into a vector as:

$$\boldsymbol{\theta} = \{\mathbf{W}_l, \mathbf{b}_l\}$$

$$\text{Where } \dim \boldsymbol{\theta} = \sum_{l=1}^L p_{l-1} p_l + \sum_{l=1}^L p_l$$

From here we have that training loss for  $\tau = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$  is:

$$\ell_{\tau}(g(\cdot \mid \boldsymbol{\theta})) := \frac{1}{n} \sum_{i=1}^n \text{Loss}(\mathbf{y}_i, g(\mathbf{x}_i \mid \boldsymbol{\theta}))$$



# Back propagation

To make notation simpler take:

$$C_i(\boldsymbol{\theta}) = \text{Loss}(\mathbf{y}_i, g(\mathbf{x}_i | \boldsymbol{\theta}))$$

Hence we have:

$$\ell_{\tau}(g(\cdot | \boldsymbol{\theta})) = \frac{1}{n} \sum_{i=1}^n C_i(\boldsymbol{\theta})$$

We now wish to find  $\nabla_{\boldsymbol{\theta}} \ell_{\tau}(g(\cdot | \boldsymbol{\theta}))$ . Denote  $\mathbf{S}'_l(\mathbf{z}_l) = \begin{bmatrix} S'(z_1) & S'(z_2) & \dots & S'(z_{p_l}) \end{bmatrix}$ . Define:

$$\mathbf{D}_l := \text{Diag}(\mathbf{S}'_l(\mathbf{z}_l))$$

# Back propagation

We then can prove that:

$$\frac{\partial C}{\partial \mathbf{W}_l} = \delta_l \mathbf{a}_{l-1}^T$$

and

$$\nabla_{\mathbf{b}_l} C = \delta_l$$

Where, we have the recursion:

$$\delta_{l-1} = \mathbf{D}_{l-1} \mathbf{W}_l^T \delta_l \implies \delta_{l-1} = \text{Diag}(\mathbf{S}'_l(\mathbf{z}_l)) \odot \mathbf{W}_l^T \delta_l$$

With base case:

$$\delta_L = \frac{\partial \mathbf{S}_L}{\partial \mathbf{z}_L} \nabla_g C$$

# Back propagation

Hence we can now find  $\nabla_{\theta} C$  as:

---

## Algorithm 2 Computing gradient of cost

---

**Input:**  $\mathbf{x} \in \mathbb{R}^{p_0}$ , corresponding  $\mathbf{y}$ ,  $\{\mathbf{W}_l \in \mathbb{R}^{p_l \times p_{l-1}}\}$ ,  $\{\mathbf{b}_l \in \mathbb{R}^{p_l}\}$ ,  $\{\mathbf{S}_l\}$

**Output:**  $\nabla_{\theta} C$

- 1:  $\{\mathbf{a}_l\}, \{\mathbf{z}_l\} \leftarrow$  Algorithm 1
  - 2:  $\delta_L \leftarrow \frac{\partial \mathbf{S}_L}{\partial \mathbf{z}_L} \nabla_g C$
  - 3:  $\mathbf{z}_0 \leftarrow \mathbf{0}$
  - 4: **for**  $l = L$  backwards to 1 **do**
  - 5:      $\nabla_{\mathbf{b}_l} C \leftarrow \delta_l$
  - 6:      $\frac{\partial C}{\partial \mathbf{W}_l} \leftarrow \delta_l \mathbf{a}_{l-1}^T$
  - 7:      $\delta_{l-1} \leftarrow \text{Diag}(\mathbf{S}'_l(\mathbf{z}_l)) \odot \mathbf{W}_l^T \delta_l$
  - 8: **end for**
  - 9: Collect  $\{\nabla_{\mathbf{b}_l} C\}$  and  $\{\frac{\partial C}{\partial \mathbf{W}_l}\}$  into  $\nabla_{\theta} C$
  - 10: **return**  $\nabla_{\theta} C$
- 

- ▷ Forward propagation
- ▷ Backward propagation

Hence we can find  $\nabla_{\theta} \ell_{\tau}(g(\cdot | \theta))$  by simply looping over Algorithm 2 for each pair in  $\tau = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$  adding up the obtained derivatives for each  $C_i$ . Now we can work on finding  $\theta^* := \operatorname{argmin}_{\theta} \ell_{\tau}(g(\cdot | \theta))$ . For brevity  $\nabla_{\theta} \ell_{\tau} := \nabla_{\theta} \ell_{\tau}(g(\cdot | \theta))$ . We take a steepest descent method of the form:

$$\theta_{t+1} = \theta_t - \alpha_t \nabla_{\theta_t} \ell_{\tau}$$

We wish to use  $\alpha_t$  to control step size: higher curvature means we should have slower rate of descent whereas low curvature means we should have higher rate of descent. Normally we would take:  $\alpha_t = (H_{\theta} \ell_{\tau})^{-1}$ . However, this is computationally expensive. Instead we take the short Barzilai-Borwein formula ( $\delta = \theta_{t+1} - \theta_t$  and  $\mathbf{s} = \nabla_{\theta_{t+1}} \ell_{\tau} - \nabla_{\theta_t} \ell_{\tau}$ ):

$$\alpha_t = \frac{\delta^T \mathbf{s}}{\|\mathbf{s}\|}$$

---

## Algorithm 3 Gradient descent training

---

**Input:**  $\tau = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ , initial guess  $\theta_0$  in terms of  $\{\mathbf{W}_l \in \mathbb{R}^{p_l \times p_{l-1}}\}$  and  $\{\mathbf{b}_l \in \mathbb{R}^{p_l}\}, \{\mathbf{S}_l\}$

**Output:**  $\operatorname{argmin}_{\theta} \ell_{\tau}(g(\cdot | \theta))$

```
1:  $t \leftarrow 1$    $\delta \leftarrow 0.1 \cdot \mathbf{1}$    $\nabla_{\theta_{t-1}} \ell_{\tau} \leftarrow \mathbf{0}$    $\alpha \leftarrow 0.1$  ▷ Initialise parameters suitably
2: while stopping condition is not met do
3:    $\nabla_{\theta_t} \ell_{\tau} \leftarrow \text{BackPropogation}$ 
4:    $\mathbf{s} \leftarrow \nabla_{\theta_t} \ell_{\tau} - \nabla_{\theta_{t-1}} \ell_{\tau}$ 
5:   if  $\delta^T \mathbf{s} > 0$  then
6:      $\alpha \leftarrow \frac{\delta^T \mathbf{s}}{\|\mathbf{s}\|}$ 
7:   else
8:      $\alpha \leftarrow \alpha \cdot 2$ 
9:   end if
10:   $\delta \leftarrow -\alpha \nabla_{\theta_t} \ell_{\tau}$    $\theta_{t+1} \leftarrow \theta_t + \delta$    $t \leftarrow t + 1$ 
11: end while
12: return  $\theta_t$ 
```

# Training

Now we look at stochastic gradient descent. Consider the case that the training set  $\tau = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ . In the case that  $n$  is very large, the training will be expensive. Now consider a discrete R.V  $K$  such that  $\forall k = 1, 2, \dots, n$

$$\mathbb{P}[K = k] = \frac{1}{n}$$

We hence have that, by definition,

$$\ell_{\tau}(g(\cdot | \boldsymbol{\theta})) := \frac{1}{n} \sum_{i=1}^n \text{Loss}(\mathbf{y}_i, g(\mathbf{x}_i | \boldsymbol{\theta})) = \mathbb{E}[\text{Loss}(\mathbf{y}_K, g(\mathbf{x}_K | \boldsymbol{\theta}))]$$

Hence, for some  $N \ll n$  i.i.d  $K_i \sim K$  we have the Monte Carlo estimate:

$$\hat{\ell}_{\tau}(g(\cdot | \boldsymbol{\theta})) := \frac{1}{N} \sum_{i=1}^N \text{Loss}(\mathbf{y}_{K_i}, g(\mathbf{x}_{K_i} | \boldsymbol{\theta}))$$

# Polynomial Regression

We consider an architecture of:

$$\begin{bmatrix} p_0 & p_1 & p_2 & p_3 \end{bmatrix} = \begin{bmatrix} 1 & 20 & 20 & 1 \end{bmatrix}$$

This results in:

Parameters between  $p_0$  and  $p_1$ :  $p_0 \times p_1 = 20$  weights and  $p_1 = 20$  biases.

Parameters between  $p_1$  and  $p_2$ :  $p_1 \times p_2 = 400$  weights and  $p_2 = 20$  biases.

Parameters between  $p_2$  and  $p_3$ :  $p_2 \times p_3 = 20$  weights and  $p_3 = 1$  bias.

Hence,  $\dim \theta = 481$ . Moreover, we can use a ReLU activation function.

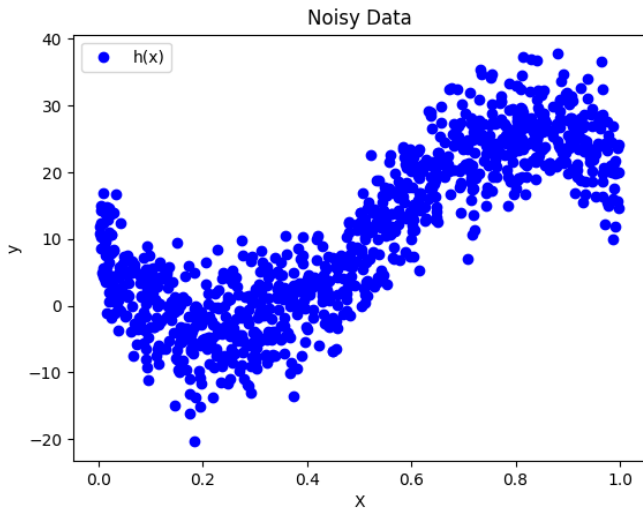
# Polynomial Regression

First we generate data:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 np.random.seed(10)
5
6 def GenerateData(NumberOfPoints : int):
7     X = np.random.uniform(0, 1, NumberOfPoints)
8     h = lambda x: 10 - 140*x + 400*(x**2) - 250*(x**3)
9     y = np.random.normal(h(X), np.sqrt(25))
10    return (X,y)
11
12 (X,y) = GenerateData(1000)
13
14 X = X.reshape(-1,1)
15 y = y.reshape(-1,1)
```



# Polynomial Regression



# Polynomial Regression

Now we can start writing the neural network. Consider the activation function:

```
1 def ReLU(z, l : int):  
2     # If last layer, dont apply ReLU  
3     if l == 3:  
4         return (z, np.ones_like(z))  
5     # Returns (S(z), S'(z)) = (z 1_{z >= 0}, 1_{z >= 0}) element-wise  
6     else:  
7         # Converted to bool we have True = 1.0 and False = 0.0.  
8         # Derivative is diagonal rep.  
9         return (np.maximum(0,z), np.array(z > 0, dtype=float))
```

# Polynomial Regression

Now we initialize the weights as:

```
1 def Initialise(p):
2     W = [None] * len(p)
3     b = [None] * len(p)
4
5     for l in range(1, len(p)):
6         # Weight matrix has to be of dimension p[l] \times p[l-1]
7         W[l] = np.random.randn(p[l], p[l-1])
8         # Bias b is a vector of size p[l]
9         b[l] = np.random.randn(p[l], 1)
10    return W, b
```

# Polynomial Regression

We can write feed-forward propagation as:

```
1 def FeedForward(x,W,b):
2     # The Jacobian is again stored as diagonal matrix here
3     a, z, DdS_dz = [None] * 4, [None] * 4, [None] * 4
4     a[0] = x.reshape(-1,1)
5     # Processing like below needs to be done from layer 1 to 3
6     for l in range(1,4):
7         z[l] = W[l] @ a[l-1] + b[l]
8         a[l], DdS_dz[l] = ReLU(z[l],1)
9     return a, z, DdS_dz
```

Now we can focus on training. Taking squared error loss we write the function as:

```
1 def Loss(y,g):
2     return (g - y)**2, 2 * (g - y)
```

Like with ReLU this also returns the derivative.

# Polynomial Regerssion

```
def BackPropagation(W,b,X,y):
    n = len(y)
    Delta = [None]*4
    dC_db, dC_dW = [None]*4, [None]*4
    LossIncurred = 0
    # For each (x,y) in the training
    # set:

    for i in range(n):
        a, z, DdS_dz =
            FeedForward(X[i,:].T, W, b)
        # Cost and its gradient at last
        # layer
        C , dC_dg = Loss(y[i], a[3])
        # Update loss incurred
        LossIncurred += C/n
        # Is dot product here and is
        # equivalent to full mmult as
        # dS_dz is diagonal
        Delta[3] = DdS_dz @ dC_dg
```

```
# Backwards propagation
for l in range(3,0,-1):
    dCi_dbl = Delta[l]
    dCi_dWl = Delta[l] @ a[l-1].T
    # Accumulate
    if dC_db[l] is None:
        dC_db[l] = dCi_dbl / n
    else:
        dC_db[l] += dCi_dbl / n
    if dC_dW[l] is None:
        dC_dW[l] = dCi_dWl / n
    else:
        dC_dW[l] += dCi_dWl / n
    # z[0] and DdS_dz[0] are None
    # (FeedForward dosent calculate)
    if l > 1:
        Delta[l-1] = DdS_dz[l-1] *
            (W[l].T @ Delta[l])

return dC_dW , dC_db , LossIncurred
```

# Polynomial Regression

Hence we now have everything to begin training. As done in the textbook, we can now do gradient descent. Here we do:

```
for epoch in range(1, num_epochs + 1):
    batch_idx = np.random.choice(n, batch_size)
    batch_X = X[batch_idx].reshape(-1, 1)
    batch_y = y[batch_idx].reshape(-1, 1)

    dc_dW, dc_db, batch_loss = BackPropagation(W, b, batch_X, batch_y)
    loss_arr.append(batch_loss.flatten()[0])

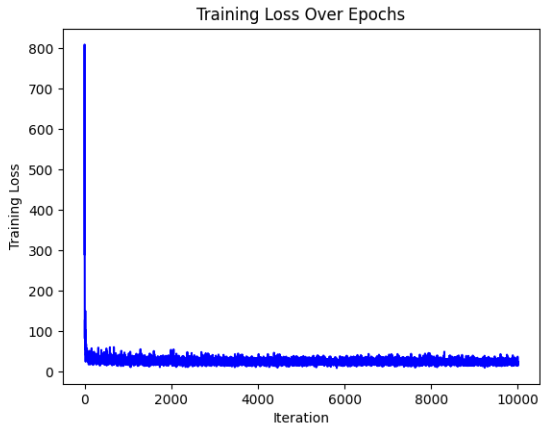
    d_beta = list2vec(dc_dW, dc_db)
    beta = list2vec(W, b)

    beta = beta - lr * d_beta

    W, b = vec2list(beta, p)
```

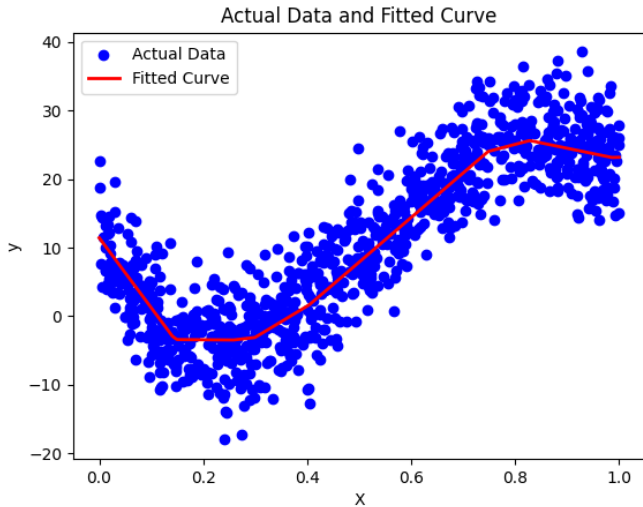
Taking batch size = 50

# Result



```
epoch | batch loss
-----
1: 291.0456945346262
1000: 29.106530603335138
2000: 27.909701127562716
3000: 19.162703429718366
4000: 24.478730987696775
5000: 20.77197762268555
6000: 17.870963762181447
7000: 21.59237202881953
8000: 18.538087981621157
9000: 17.72078710219764
10000: 22.82676949209595
```

# Result



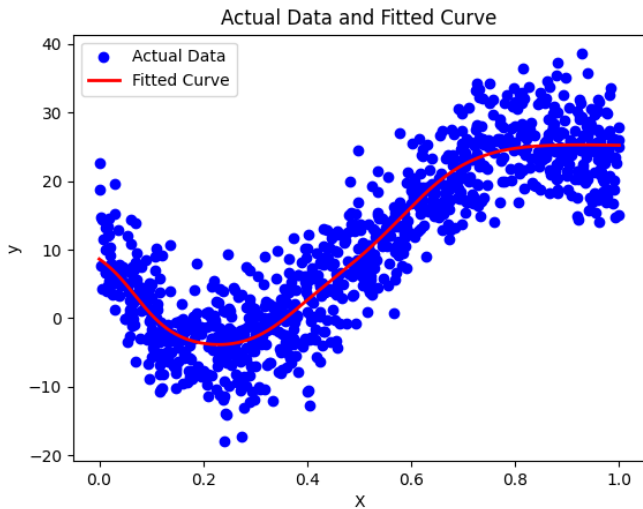


## Bonus: Using sigmoid activation

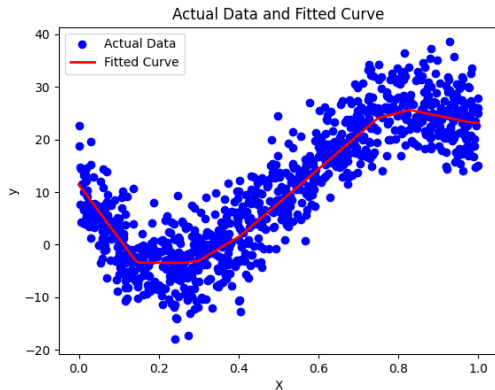
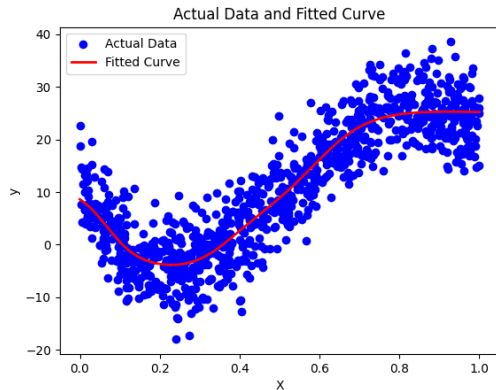
To do so you can just use this function instead:

```
1 def Sigmoid(z, l : int):  
2     if l == 3:  
3         return (z, np.ones_like(z))  
4     else:  
5         S = lambda z: 1/(1+np.exp(-z))  
6         Sprime = lambda z: (np.exp(z))/((1+np.exp(z))**2)  
7         return (S(z), Sprime(z))
```

## Bonus: Using sigmoid activation



# Bonus: Using sigmoid activation

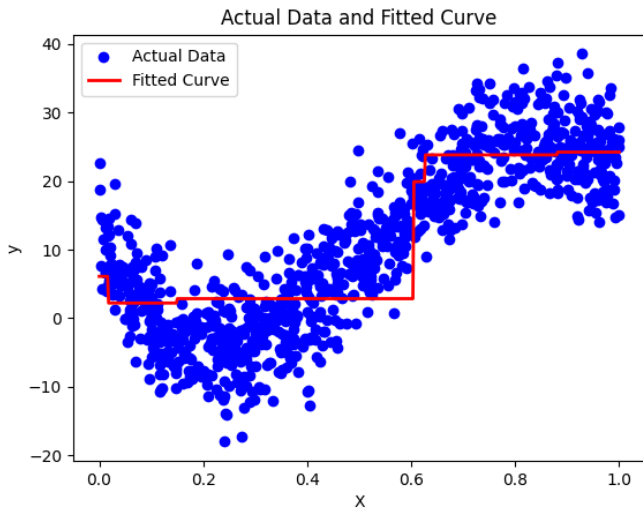


## Bonus: Using indicator activation

To do so you can just use this function instead:

```
1 def Indicator(z, l : int):  
2     if l == 3:  
3         return (z, np.ones_like(z))  
4     else:  
5         return (np.array(z > 0, dtype=float), np.zeros_like(z))
```

## Bonus: Using indicator activation



# Bonus: Using indicator activation

