# Lecture 05
# Deep Learning: Basics

February 20, 2025

Data Science Using Python

**Jaeung Sim**

Assistant Professor

School of Business, University of Connecticut

# Agenda

- **Python Hands-on for Concepts in Lecture 04**

- **Q&A for Data Collection Methods**

- **Deep Learning: Basics**

- **Announcements**

# Python Hands-on for Concepts in Lecture 04

**Regression and Predictive Analysis**

# Q&A for Data Collection Methods

**Beautiful Soup and Selenium**

# Deep Learning: Basics

**Basic Conceptual Background**

# Neural Network

- What is a **Neural Network**?

  - A neural network is a type of machine learning algorithm that is modeled after the structure and function of the human brain.

  - It is composed of interconnected **nodes** or artificial **neurons** that are organized in layers.



(a) Biological neuron

(b) Artificial neuron

# Neural Network

- What is **Deep Learning**?
  - Deep learning uses artificial neural networks with **many layers** (thus, "**deep**") to model and solve complex problems.
  - It is based on the idea of learning hierarchical representations of data, where each layer of the network learns to represent increasingly abstract features of the input data.
  - Deep learning algorithms require large amounts of data and computational power.



**Deep Neural Network**

input layer   hidden layer 1   hidden layer 2   hidden layer 3   output layer

Figure 12.2 Deep network architecture with multiple layers.

# Neural Network

- Examples of Deep Learning
  - **Convolutional Neural Network (CNN)**
    - Image classification, object detection, video action recognition

# Neural Network

- Examples of Deep Learning

    - **Recurrent Neural Network (RNN)**

        - Text classification and generation

        - Time series forecasting

# Neural Network

- **Structure**
  - **Nodes**
  - **Layers**
  - **Inputs ($X$)**
  - **Weights ($w$)**
  - **Bias ($b$)**
  - **Summation Function ($Z$)**
  - **Activation Function ($f(Z)$)**

**Hidden layer**

$f(Z_1)$

**Node**

**Input layer**

**Output layer**

$Z_1 = w_1 X + b_1$

$w_3 f(Z_1)$

**X**

$w_3 f(Z_1) + w_4 f(Z_2) + b_3$

**Y**

$Z_2 = w_2 X + b_2$

$w_4 f(Z_1)$

10

# Neural Network

- **Structure**
  - **Inputs:** Inputs are the set of values for which we need to predict an output value. They can be viewed as features or attributes in a dataset.
  - **Weights:** Weights are the real values that are attached with each input/feature and they convey the importance of that corresponding feature in predicting the final output.
  - **Bias:** Bias is used for shifting the activation function towards left or right, you can compare this to y-intercept in the line equation.
  - **Summation Function:** The work of the summation function is to bind the weights and inputs together and calculate their sum.
  - **Activation Function:** It is used to introduce non-linearity in the model.

# Neural Network

- **Estimation Steps**

  - Initialize the weights and biases with some random values.

  - Calculate the output which is to predict the values and estimate the loss.

    - Refer to: ***Activation function***

  - Adjust the weights and biases such that the loss will be minimized.

    - Refer to: ***Gradient descent***, ***backpropagation***

# Basic Example

- **Bag Price and Purchase Decision**

# Basic Example

- **Bag Price and Purchase Decision**

# Basic Example

- **Estimation Steps**
  - Let's start with some given parameters.
  - Calculate the output and loss.
  - Adjust weight/loss.

**Softplus:** $f(Z_1) = log(1 + e^{Z_1})$

$\times (-34.4)$

$+ 2.14$

$\times (-1.30)$

**Price**

Input

$\times (-2.52)$

$+ 1.29$

$\times 2.28$

Weighted Sum

$+ (-0.58)$

**Purchase**

Output

# Activation Function

- Also known as a **transfer function**

- Mapping a linear function to another function

**Linear Activation Function**



**Doesn't help with the complexity
in neural networks**

**Nonlinear Activation Functions**



**Makes it easy for the model to
adapt with variety of data**

# Activation Function

- **Sigmoid function**
    - Exists **between 0 and 1**.
    - Especially used for models to predict the probability as an output.
    - The function is **differentiable**.
    - The function is **monotonic**, but the function's derivative is not.
    - The logistic sigmoid function can cause a neural network to get stuck at the training time.
    - The softmax function is a more generalized logistic activation function which is used for multiclass classification.

$$\phi(z) = \frac{1}{1 + e^{-z}}$$

# Activation Function

- **Hyperbolic Tangent (tanh)**
  - Exists **between -1 and 1**.
  - The tanh function is mainly used classification between two classes.
  - The negative inputs are mapped strongly negative, and the zero inputs are mapped near zero in the tanh graph.
  - The function is **differentiable**.
  - The function is **monotonic** while its derivative is not monotonic.

# Activation Function

- **Rectified Linear Unit (ReLU)**

  - $R(x) = \max(0, x)$

  - Exists from 0 to infinity

  - **The most used activation function**

  - Half rectified (from bottom)

  - The function and its derivative both are monotonic.

  - All the negative values become zero immediately which decreases the ability of the model to fit or train from the data properly.



sigmoid

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



ReLU

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

# Activation Function

- **Softplus**
  - $f(x) = \log(1 + \exp(x))$
  - A smooth approximation to the ReLU activation function.
  - As $x \to -\infty$, it becomes identical to
    $\text{sigmoid}(x) = \frac{1}{1 + e^{-x}}$

# Activation Function

- **Application to the Basic Example**



**Softplus:** $f(Z_1) = log(1 + e^{Z_1})$

Price
Input

$\times (-34.4)$

$+ 2.14$

$\times (-1.30)$

$\times (-2.52)$

$+ 1.29$

$\times 2.28$

Weighted Sum

$+ (-0.58)$

**Purchase**
Output

# Activation Function

- **Application to the Basic Example**

$$f(2.14) = \log(1 + e^{2.14}) = 2.25$$

**Softplus:** $f(Z_1) = log(1 + e^{Z_1})$



$Z_1(0) = (-34.4) \times 0 + 2.14 = 2.14$

$+\,2.14$

$\times\,(-1.30)$

$\times\,(-34.4)$

**Price**

$0$

**Purchase**

$-1.30 \times 2.25 + 2.28 \times 1.53 = $ Weighted Sum

$+\,(-0.58)$

-0.011

$\times\,(-2.52)$

$\times\,2.28$

$+\,1.29$

$Z_2(0) = (-2.52) \times 0 + 1.29 = 1.29$

$f(1.29) = \log(1 + e^{1.29}) = 1.53$

# Activation Function

- **Application to the Basic Example**

$$f(-15.06) = \log(1 + e^{-15.06}) = 2.88 \times 10^{-7}$$

**Softplus:** $f(Z_1) = \log(1 + e^{Z_1})$

$$Z_1(0) = (-34.4) \times 0.5 + 2.14 = -15.06$$

$+ 2.14$

$\times (-1.30)$

$\times (-34.4)$

**Price**

0.5

$-1.30 \times 2.88 \times 10^{-7} + 2.28 \times 0.71 =$ Weighted Sum

$+ (-0.58)$

**Purchase**

1.035

$\times (-2.52)$

$\times 2.28$

$$Z_2(0) = (-2.52) \times 0.5 + 1.29 = 0.03$$

$+ 1.29$

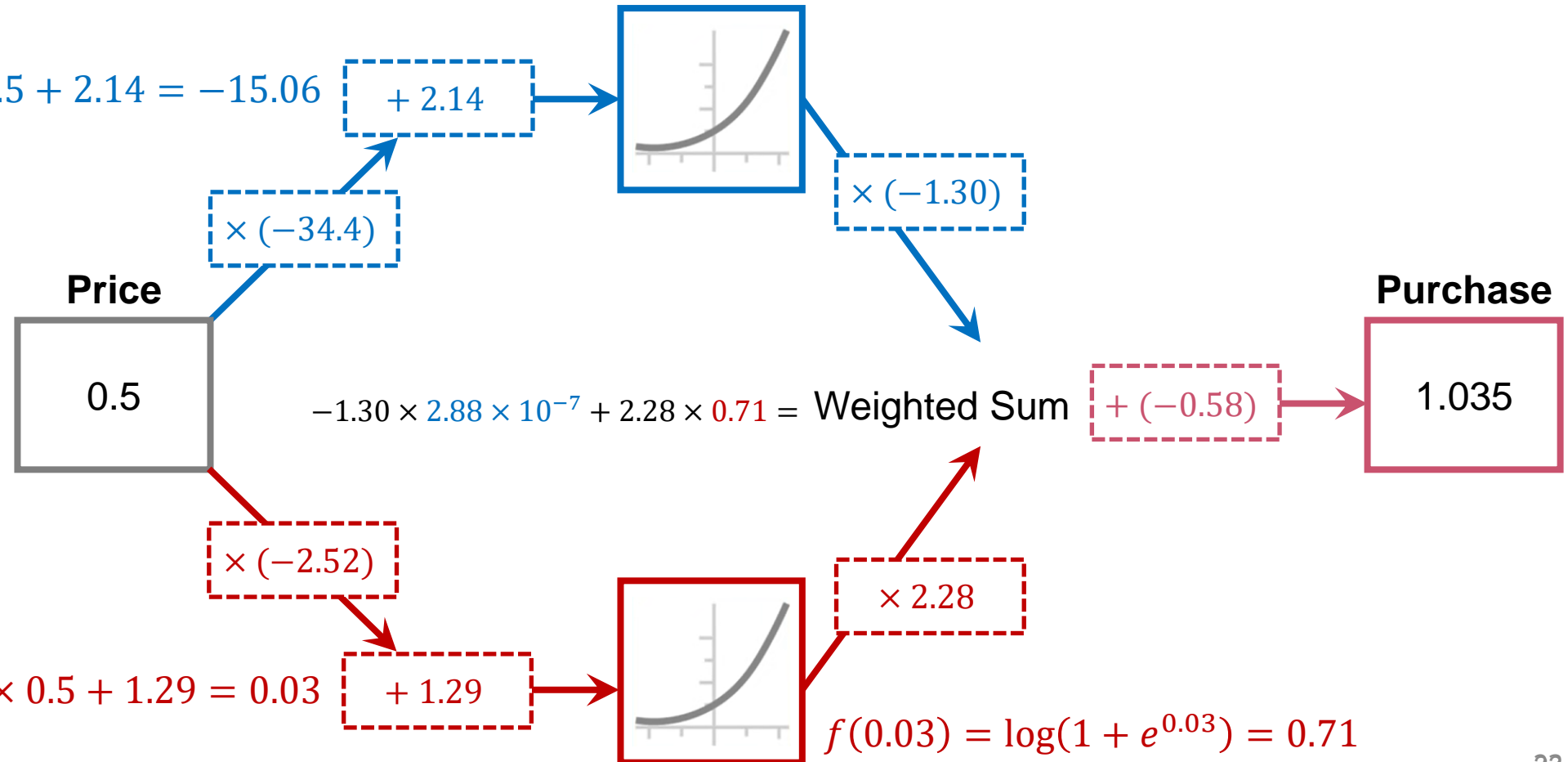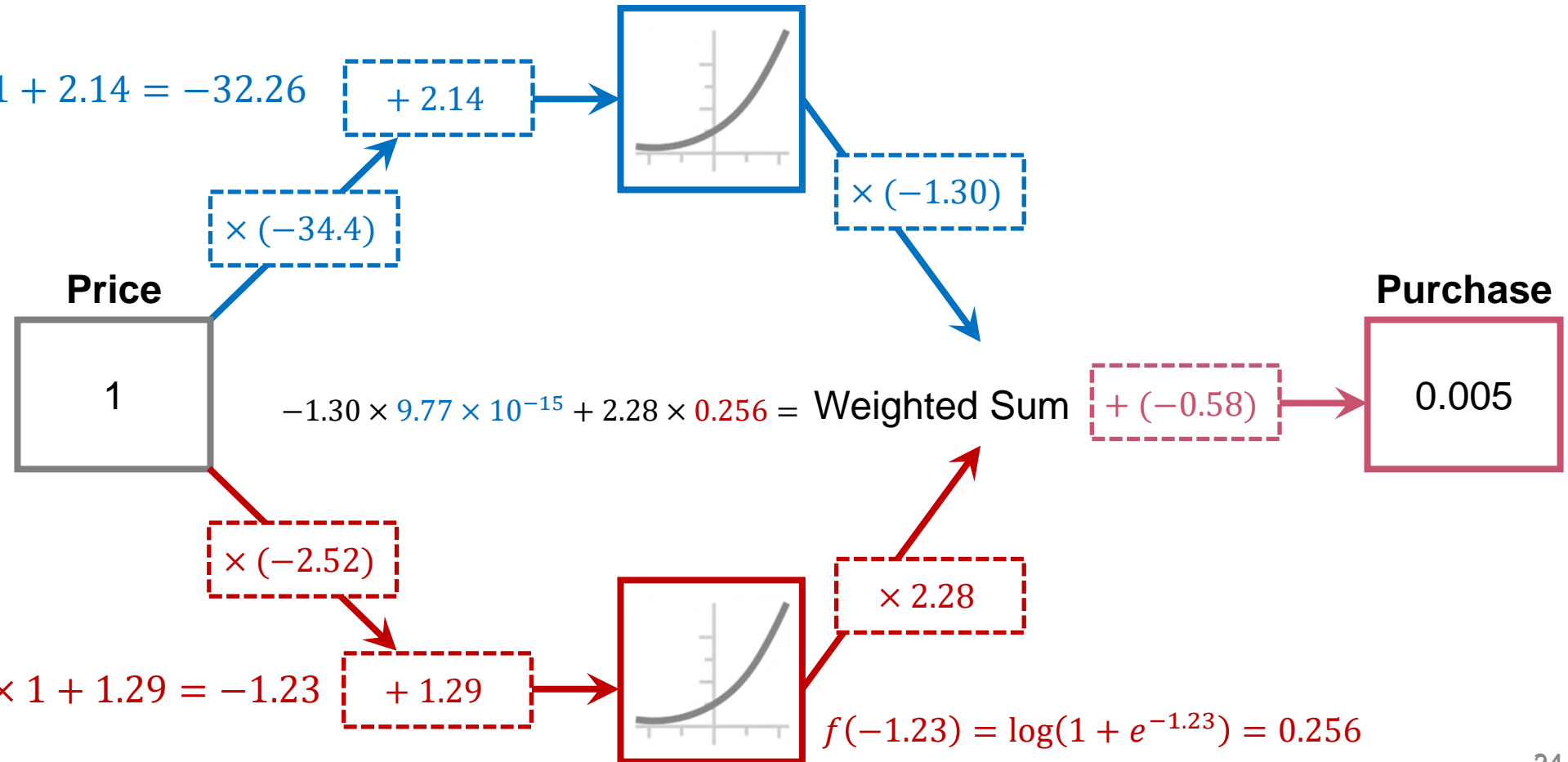$$f(0.03) = \log(1 + e^{0.03}) = 0.71$$

# Activation Function

- **Application to the Basic Example**

$$f(-32.26) = \log(1 + e^{-15.06}) = 9.77 \times 10^{-15}$$

**Softplus:** $f(Z_1) = log(1 + e^{Z_1})$

$$Z_1(0) = (-34.4) \times 1 + 2.14 = -32.26$$

$+ 2.14$

$\times (-34.4)$

$\times (-1.30)$

**Price**

**Purchase**

$1$

$-1.30 \times 9.77 \times 10^{-15} + 2.28 \times 0.256 =$ Weighted Sum

$+ (-0.58)$

$0.005$

$\times (-2.52)$

$\times 2.28$

$$Z_2(0) = (-2.52) \times 1 + 1.29 = -1.23$$

$+ 1.29$

$$f(-1.23) = \log(1 + e^{-1.23}) = 0.256$$

**(1) Linear function**

**(3) Weighting**

Softplus: $f(Z_1) = log(1 + e^{Z_1})$

$+ 2.14$

$\times (-34.4)$

$\times (-1.30)$

**Price**

X

Weighted Sum $+ (-0.58)$

**Purchase**

Y

**(2) Activate Function**

**Price**

X

**Purchase**

Y

$\times (-2.52)$

$+ 1.29$

Weighted Sum $+ (-0.58)$

$\times 2.28$

**(1) Linear function**

**(2) Activate Function**

**(3) Weighting**

# (1) Linear function



# (2) Activate Function



# (3) Weighting



# (4) Adding Bias

# Chain Rule

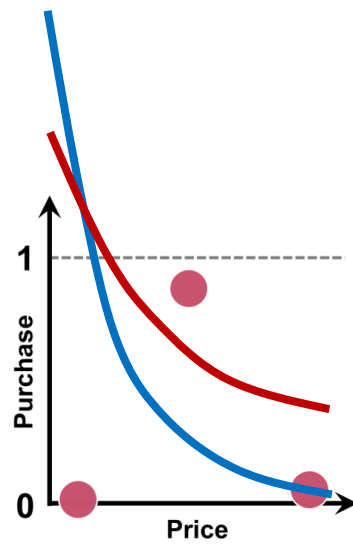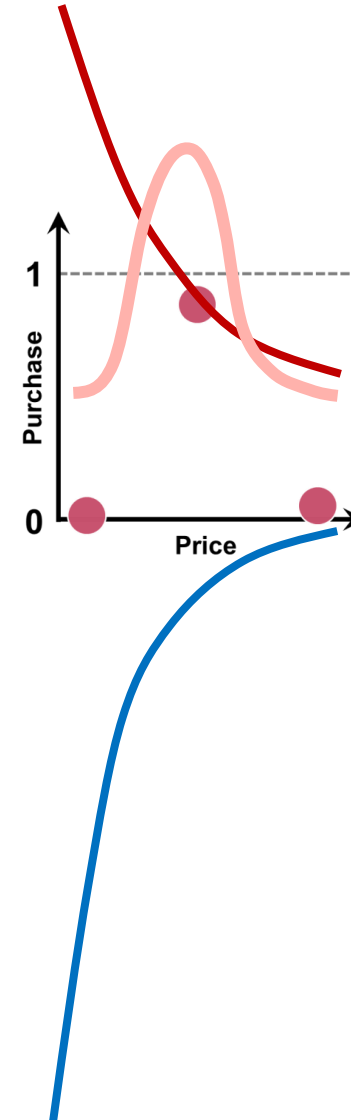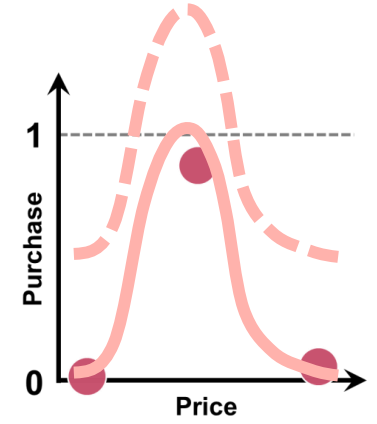- After calculating the output which is to predict the values and estimate the loss, we should adjust the weights and biases to minimize the loss.

- In doing so, we apply **backpropagation** based on **gradient descent**.

- To understand gradient descent, you need the **chain rule**:

**Lagrange's notation**

For function $h(x) = f(g(x))$,

$$h'(x) = f'\big(g(x)\big) \cdot g'(x)$$

**Leibniz's notation**

If a variable $z$ depends on the variable $y$, which itself depends on the variable $x$,

$$\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$$

# Chain Rule

- **Example: Chicken Price and Fried Chicken Sales**

# Chain Rule

- **Example: Chicken Price and Fried Chicken Sales**





We can rewrite the relationship as follows:

$$Fried\ Chicken\ Sales = f(Fried\ Chicken\ Price)$$

$$Fried\ Chicken\ Price = g(Chicken\ Price)$$

$$Fried\ Chicken\ Sales = f(g(Chicken\ Price))$$

# Chain Rule

- **Example: Chicken Price and Fried Chicken Sales**

$$\frac{d \ Fried \ Chicken \ Sales}{d \ Chicken \ Price} = \frac{d \ Fried \ Chicken \ Sales}{d \ Fried \ Chicken \ Price} \cdot \frac{d \ Fried \ Chicken \ Price}{d \ Chicken \ Price}$$

Demand-side response

*"Are consumers price-sensitive?"*

Supply-side response

*"Are sellers margin-sensitive?"*

# Gradient Descent

- A first-order iterative optimization for finding a local minimum of a differentiable function



- **Step 1:** Take the derivative of the **loss function** for each parameter in it.

- **Step 2:** Pick random values for the parameters.

- **Step 3:** Plug the parameter values into the derivatives (i.e., the **gradient**).

- **Step 4:** Calculate the step sizes as
    - **Step Size** = **Slope** x **Learning Rate**

- **Step 5:** Calculate the new parameters as
    - **New** = **Old** − **Step Size**
    - Going back to **Step 3**

# Gradient Descent

- **Loss Function**
  - **Loss** is the penalty for a bad prediction, or a number indicating how bad the model's prediction was on a single example.
  - **Regression losses:**
    - Mean squared error (MSE)/Quadratic loss/L2 loss
      - $MSE = \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{n}$
    - Mean absolute error (MAE)/L1 loss
      - $MAE = \frac{\sum_{i=1}^{n}|y_i - \hat{y}_i|}{n}$
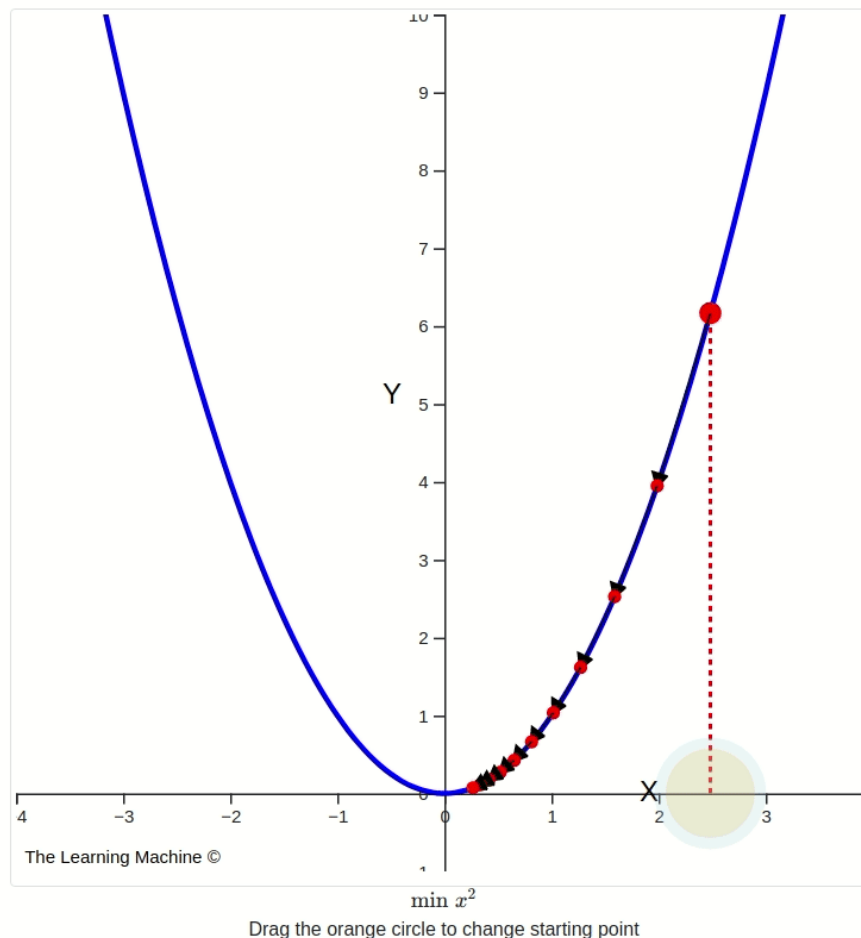  - **Classification losses:**
    - Cross entropy loss/negative log likelihood
      - $CrossEntropyLoss = -(y_i \log(\hat{y}_i) + (1 - y_i)\log(1 - \hat{y}_i))$

# Gradient Descent

- **Learning Rates**
  - The **learning rate** ($\alpha$) is a **hyperparameter** used to govern the pace at which an algorithm updates or learns the values of a parameter estimate.
  - A desirable learning rate is **low enough** for the network to converge on something useful while yet being **high enough** to train in a reasonable length of time.
  - Smaller learning rates necessitate more training epochs because of the fewer changes.
  - Larger learning rates result in faster changes and a suboptimal final set of weights frequently.
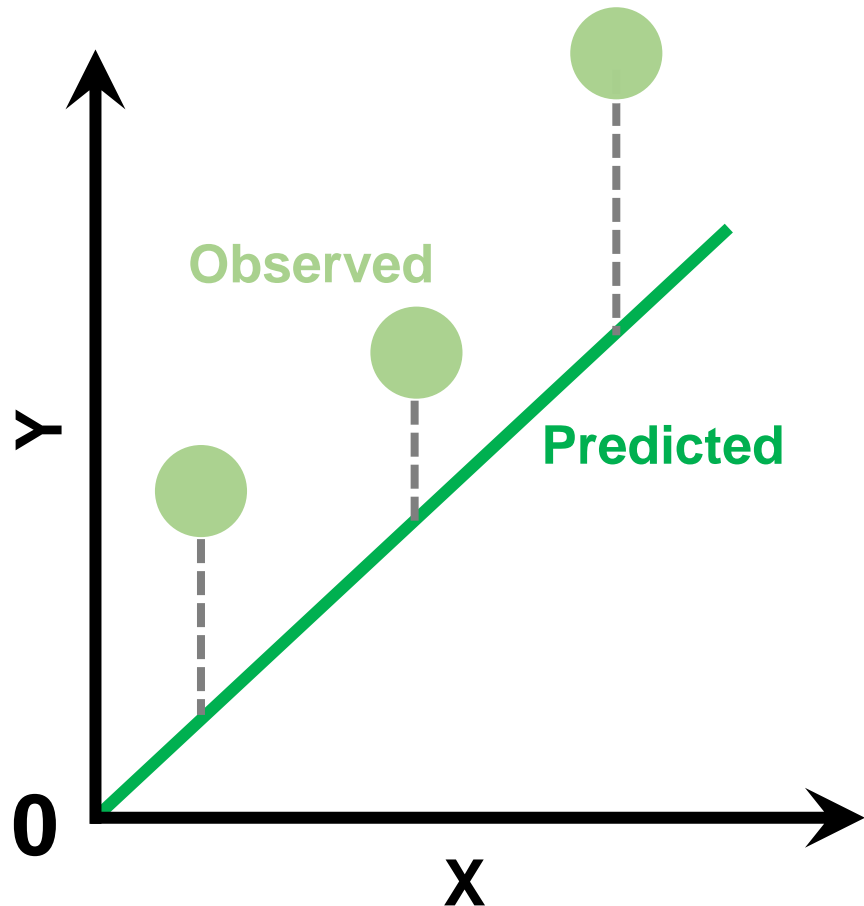


Learning rate $\alpha$

0.1

**Insights**

step 1: x = 1.99,

step 2: x = 1.59,

step 3: x = 1.27,

step 4: x = 1.02,

step 5: x = 0.814,

step 6: x = 0.651,

step 7: x = 0.521,

step 8: x = 0.417,

step 9: x = 0.333,

step 10: x = 0.267,

The Learning Machine ©

$\min x^2$

Drag the orange circle to change starting point

34

# Gradient Descent

- **A Simple Example**



**Sum of squared residuals (SSR)**

$$= \Sigma \; (\textbf{Observed} - (\textbf{Predicted}))^2$$

$$= (Y_1 - (\text{intercept} + \text{slope x } X_1))^2$$
$$+ (Y_2 - (\text{intercept} + \text{slope x } X_2))^2$$
$$+ (Y_3 - (\text{intercept} + \text{slope x } X_3))^2$$

$$= (1.4 - (\text{intercept} + 0.64 \text{ x } 0.5))^2$$
$$+ (1.9 - (\text{intercept} + 0.64 \text{ x } 2.3))^2$$
$$+ (3.2 - (\text{intercept} + 0.64 \text{ x } 2.9))^2$$

# Gradient Descent

- **A Simple Example**
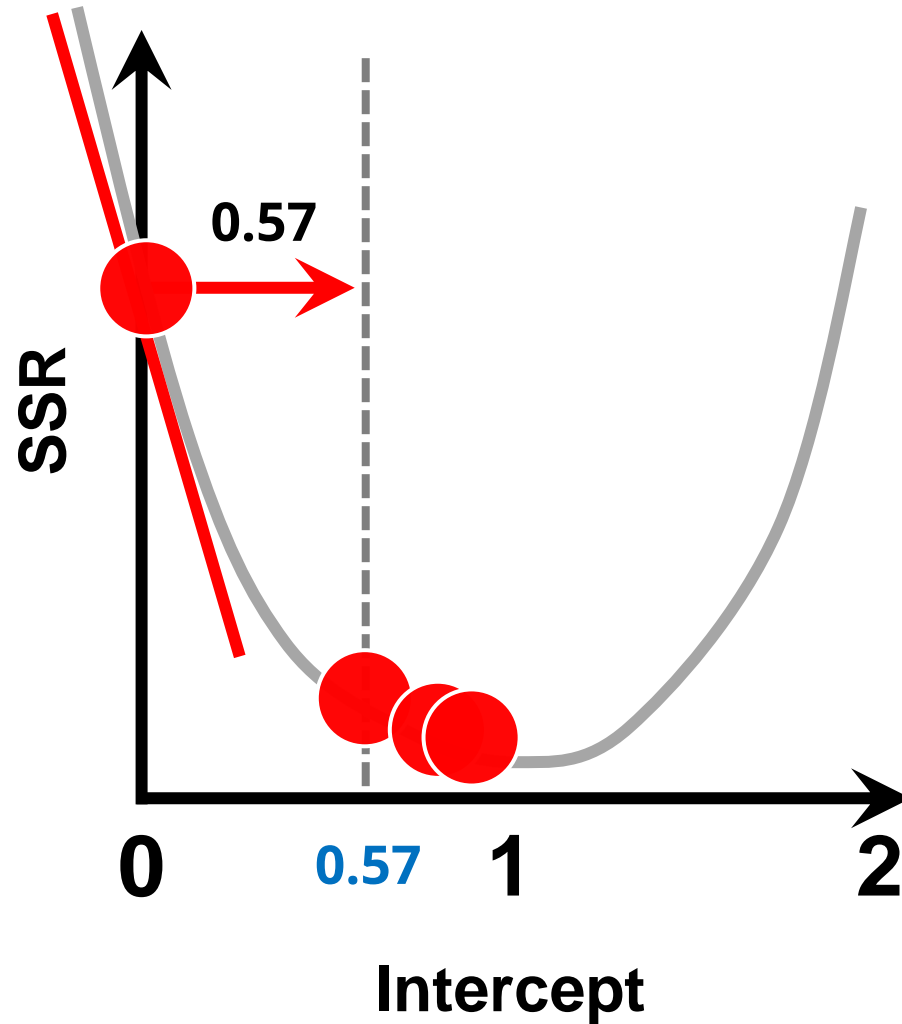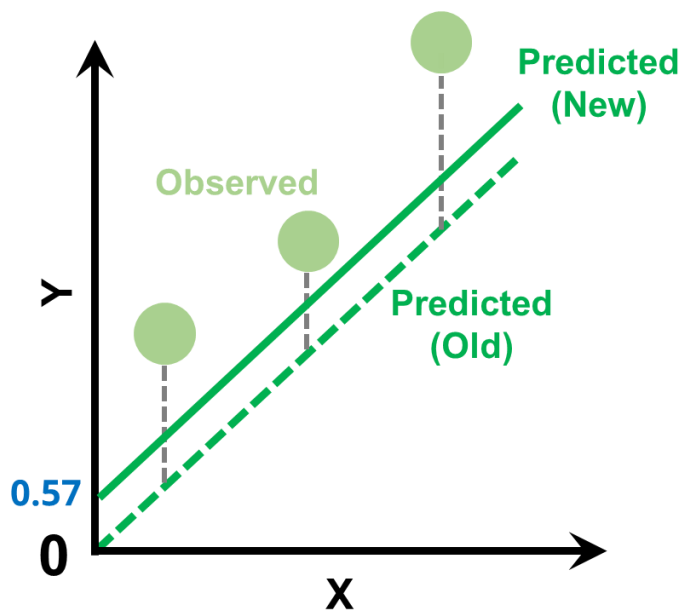
**Chain Rule:** $\dfrac{dz}{dx} = \dfrac{dz}{dy} \cdot \dfrac{dy}{dx}$

$$\frac{\partial \, SSR}{\partial \, intercept} = \frac{\partial}{\partial \, intercept} \, (1.4 - (\text{intercept} + 0.64 \text{ x } 0.5))^2 \qquad = \quad 2 \text{ x } (1.4 - (\text{intercept} + 0.64 \text{ x } 0.5)) \text{ x } (-1)$$

$$+ \frac{\partial}{\partial \, intercept} \, (1.9 - (\text{intercept} + 0.64 \text{ x } 2.3))^2 \qquad + \quad 2 \text{ x } (1.9 - (\text{intercept} + 0.64 \text{ x } 2.3)) \text{ x } (-1)$$

$$+ \frac{\partial}{\partial \, intercept} \, (3.2 - (\text{intercept} + 0.64 \text{ x } 2.9))^2 \qquad + \quad 2 \text{ x } (3.2 - (\text{intercept} + 0.64 \text{ x } 2.9)) \text{ x } (-1)$$

$$= \quad -5.7$$

36

# Gradient Descent

- **A Simple Example**

  - **Step Size = Slope x Learning Rate**

    $= -5.7 \times 0.1$

    $= -0.57$

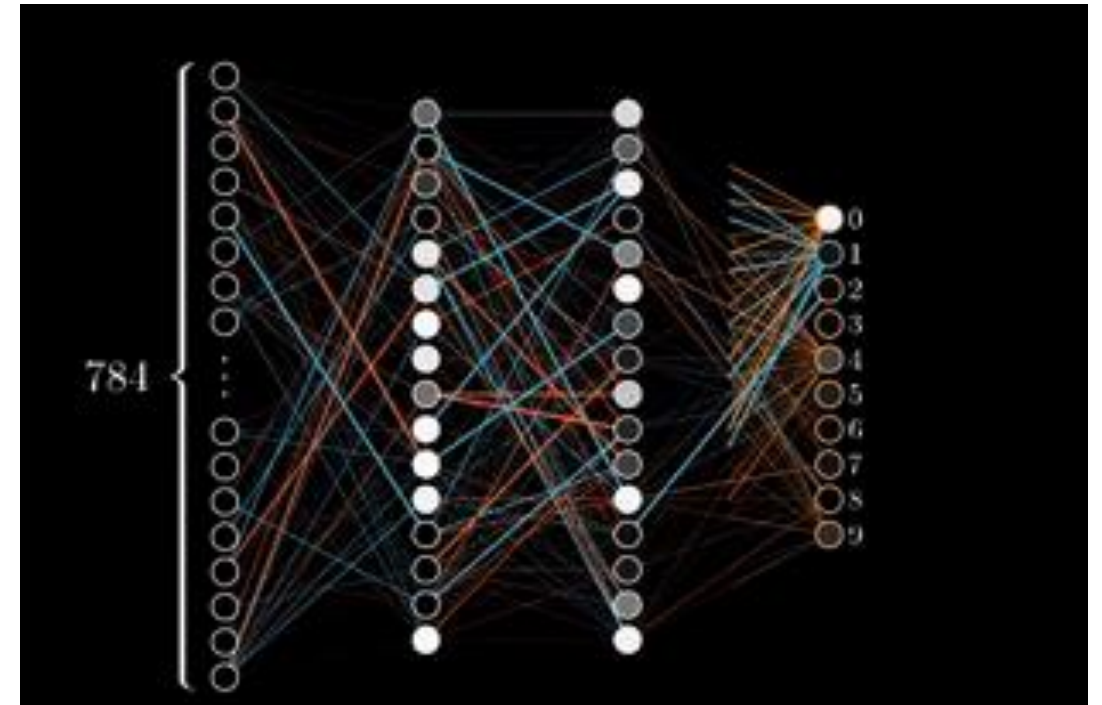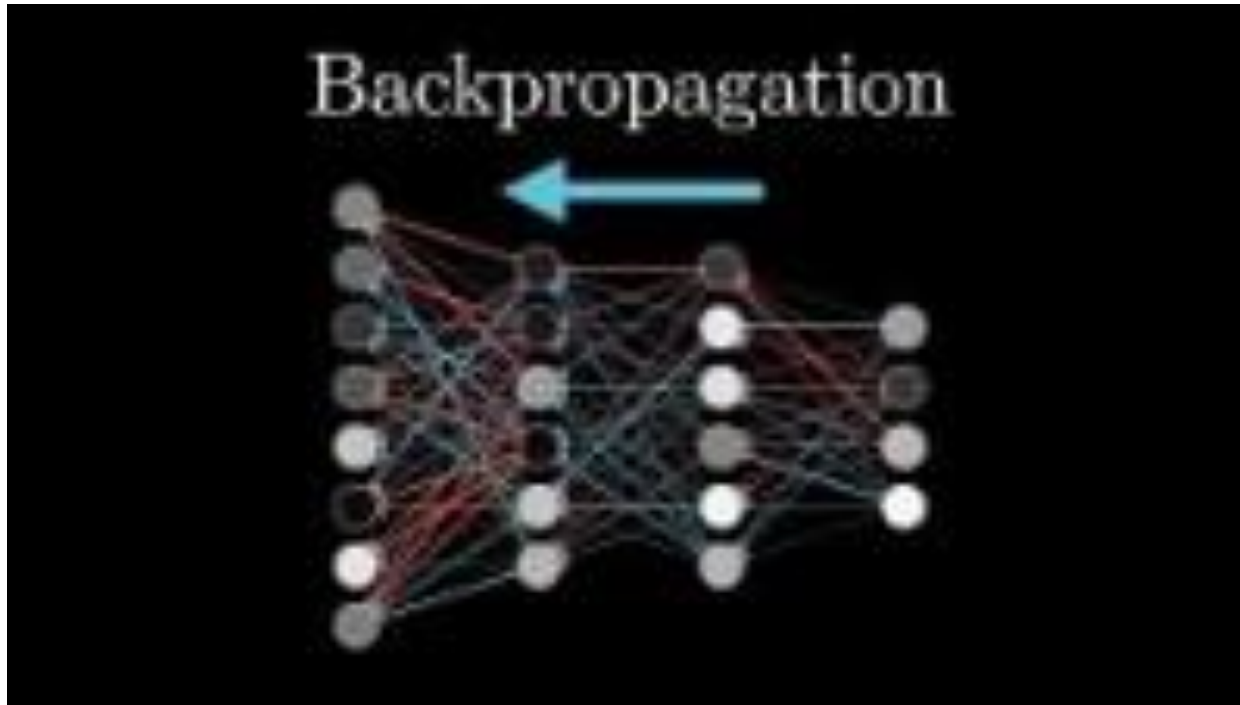  - **New Intercept = 0 – (– 0.57) = 0.57**

# Backpropagation

- Backpropagation works by **propagating the error backwards through the network**, beginning at the output layer and progressing backwards towards the input layer.

- The **gradient of the loss function** with respect to the weights is calculated using the **chain rule**, which allows the error at the output layer to be backpropagated to the hidden layers.

- The weights are then updated using the **gradient descent** algorithm, which adjusts the weights in the direction that minimizes the loss function.

- This process is repeated over many iterations until the network converges to a set of weights that minimize the loss function.

# Backpropagation

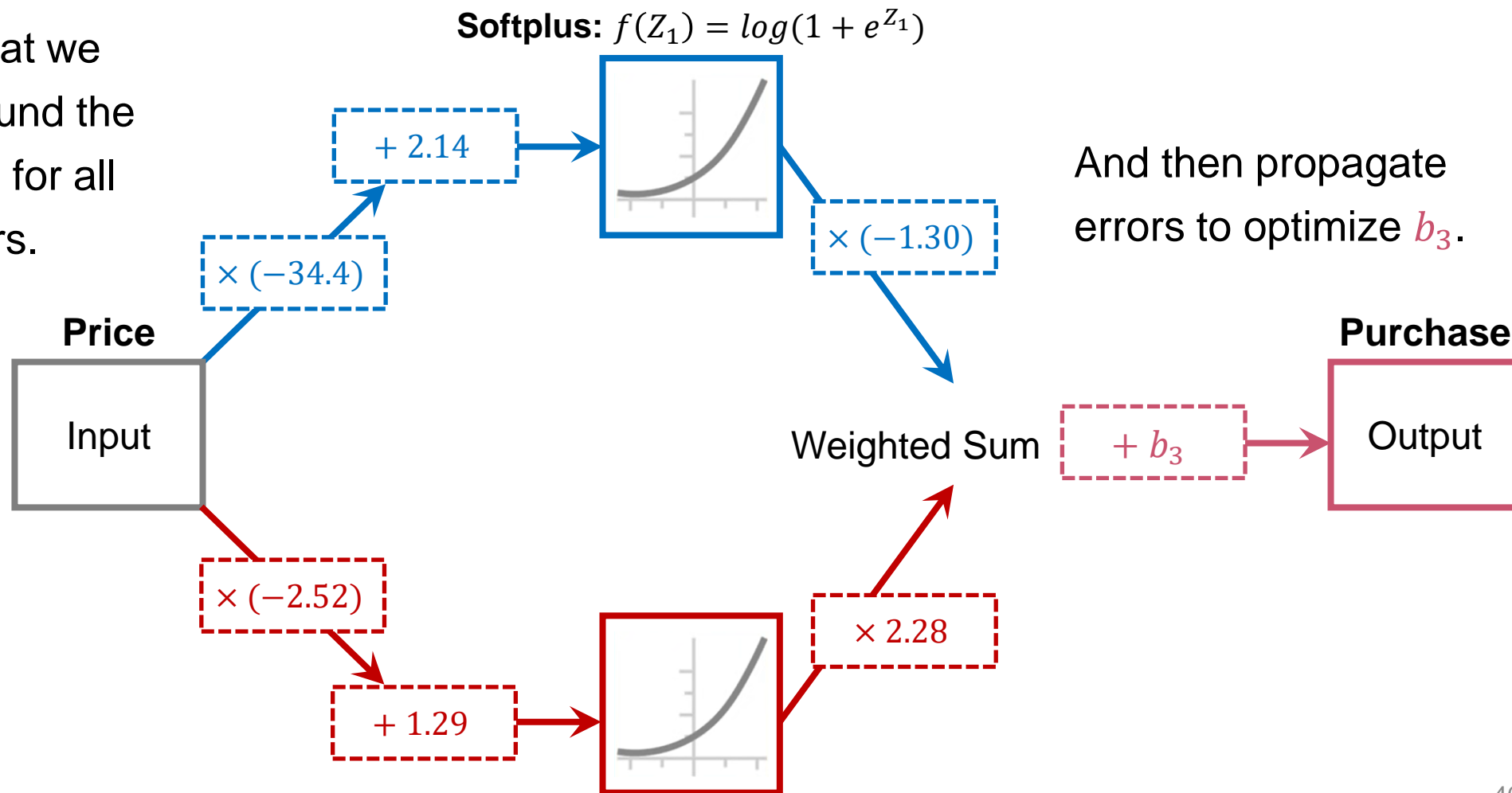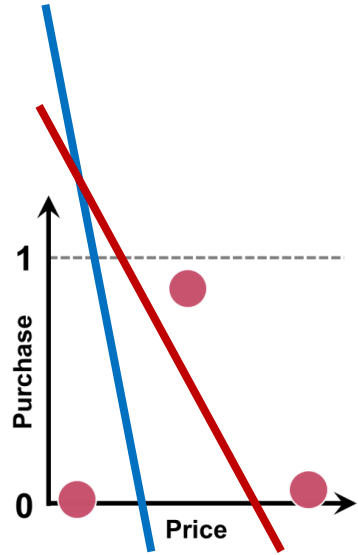- **Awesome Visualization of Backpropagation**



**Source:** https://www.youtube.com/watch?v=Ilg3gGewQ5U

# Backpropagation

- **Application to the Basic Example**

Let's assume that we have already found the optimum values for all other parameters.

**Softplus:** $f(Z_1) = log(1 + e^{Z_1})$



$+ 2.14$

$\times (-34.4)$

$\times (-1.30)$

**Price**

Input

And then propagate errors to optimize $b_3$.

$\times (-2.52)$

$+ 1.29$

$\times 2.28$

Weighted Sum
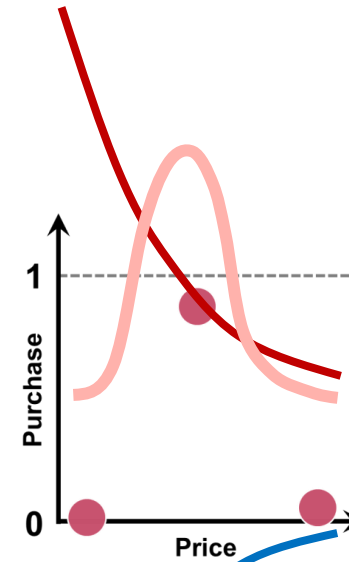
$+ b_3$

**Purchase**

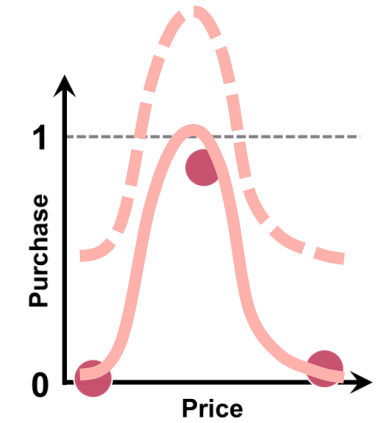Output

# (1) Linear function

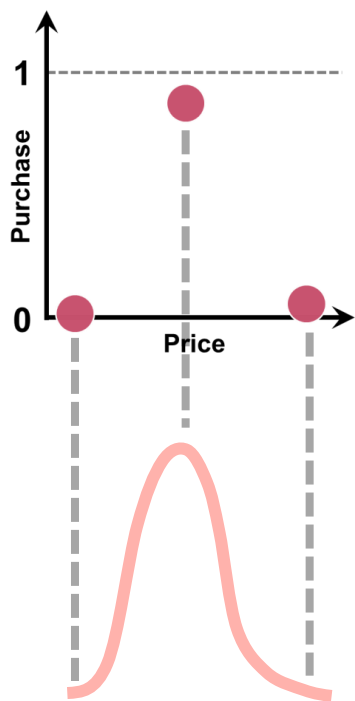

# (2) Activate Function



# (3) Weighting



# (4) Adding Bias



In other words, we propagate errors to optimize the intercept of the plot $b_3$.

# Backpropagation

- **Application to the Basic Example**



**Sum of squared residuals (SSR)**

$$= \Sigma \, (\mathbf{Observed} - (\mathbf{Predicted}))^2$$

$$= (Y_1 - (\text{intercept} + f(X_1)))^2$$
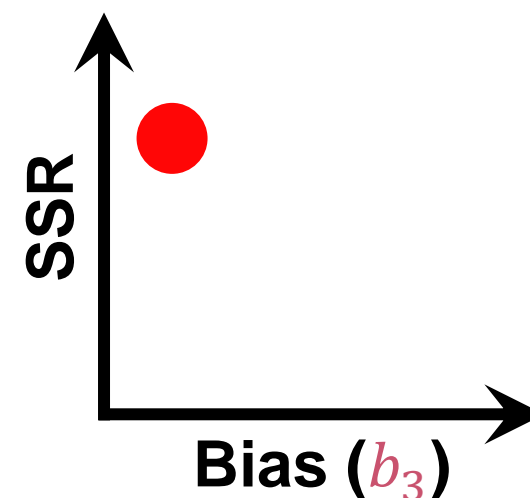$$+ (Y_2 - (\text{intercept} + f(X_2)))^2$$
$$+ (Y_3 - (\text{intercept} + f(X_3)))^2$$
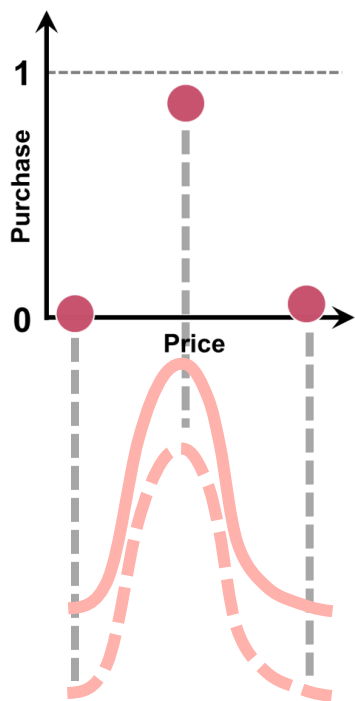
$$= (0 - (-2.6))^2$$
$$+ (1 - (-1.61))^2$$
$$+ (0 - (-2.61))^2$$

$$= \mathbf{20.38}$$

# Backpropagation

- **Application to the Basic Example**



**Sum of squared residuals (SSR)**

$$= \Sigma \ (\mathbf{Observed} - (\mathbf{Predicted}))^2$$

$$= (Y_1 - (\text{intercept} + f(X_1)))^2$$
$$+ (Y_2 - (\text{intercept} + f(X_2)))^2$$
$$+ (Y_3 - (\text{intercept} + f(X_3)))^2$$

$$= (0 - (-1.6))^2$$
$$+ (1 - (-0.61))^2$$
$$+ (0 - (-1.61))^2$$

$$= \mathbf{7.74}$$

# Backpropagation

- **Application to the Basic Example**



**Sum of squared residuals (SSR)**

$$= \Sigma \; (\textbf{Observed} - (\textbf{Predicted}))^2$$

$$= (Y_1 - (\text{intercept} + f(X_1)))^2$$
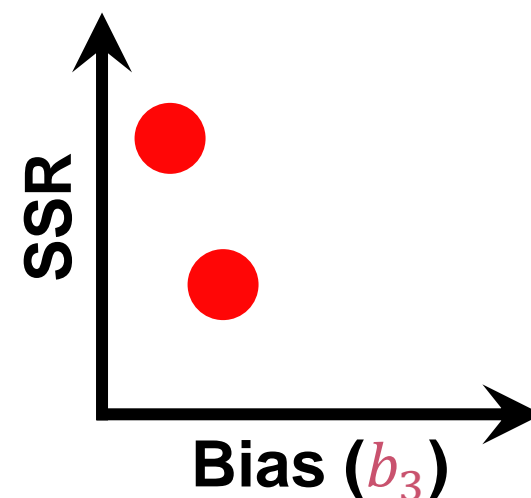$$+ (Y_2 - (\text{intercept} + f(X_2)))^2$$
$$+ (Y_3 - (\text{intercept} + f(X_3)))^2$$
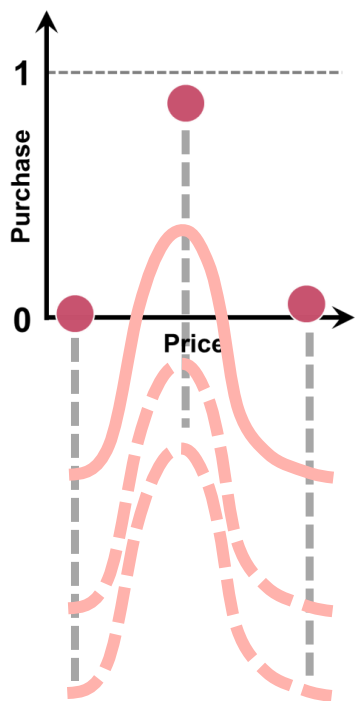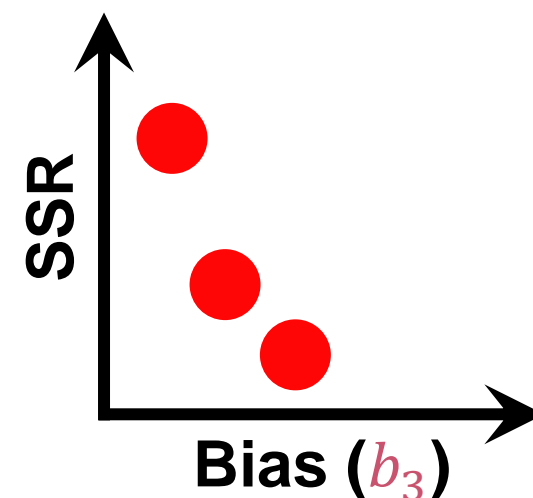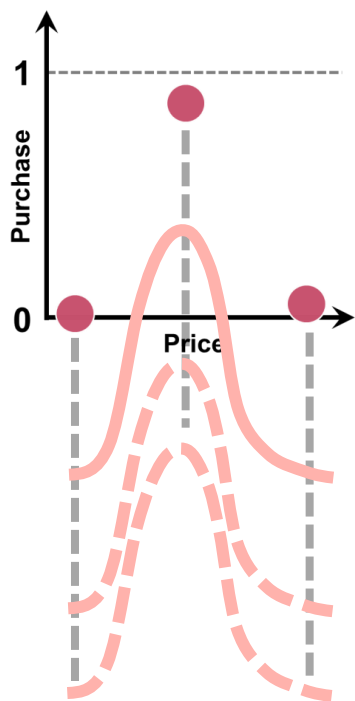
$$= (0 - (-0.6))^2$$
$$+ (1 - (0.61))^2$$
$$+ (0 - (-0.61))^2$$

$$= \textbf{0.88}$$

# Backpropagation

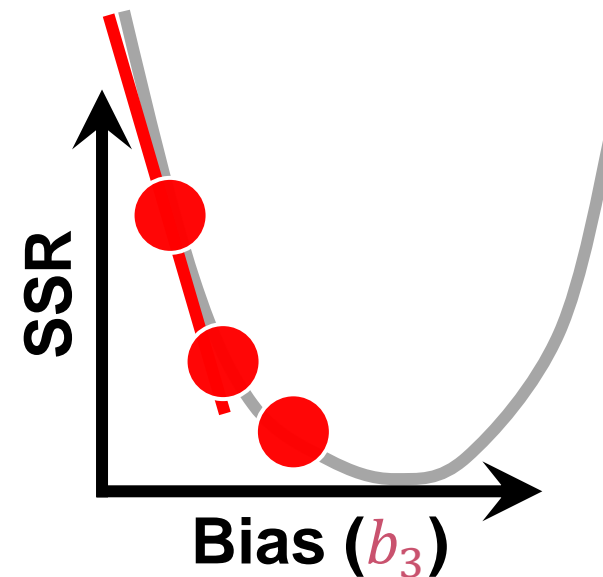- **Application to the Basic Example**



$$\frac{\partial\ SSR}{\partial\ b_3} = \frac{\partial\ SSR}{\partial\ Predicted} \cdot \frac{\partial\ Predicted}{\partial\ b_3}$$

$$= \frac{\partial}{\partial\ Predicted} \cdot \sum_{i=1}^{n=3}(Observed_i - Predicted_i)^2$$

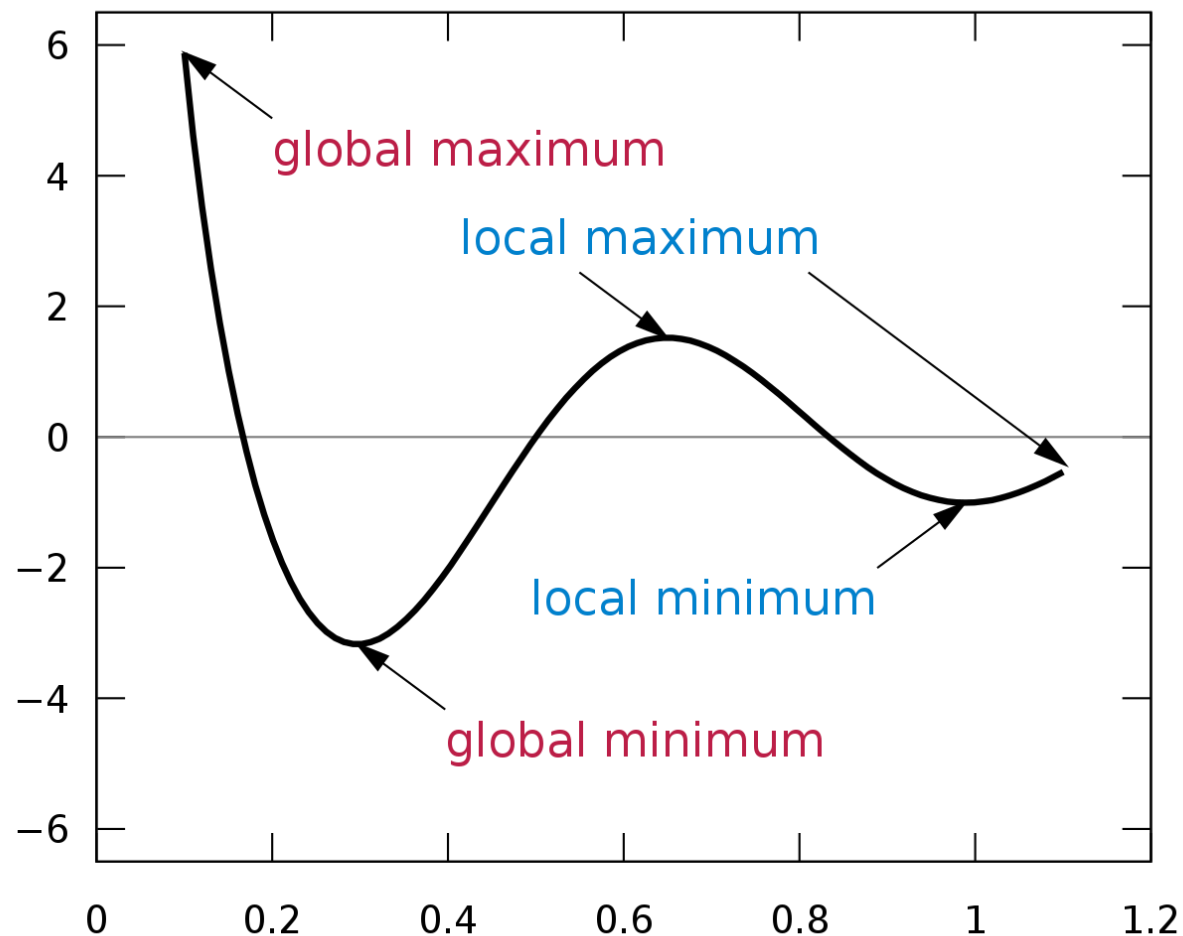$$= 2 \cdot \sum_{i=1}^{n=3}(Observed_i - Predicted_i) \cdot (-1)$$

$$\text{Step Size} = \frac{\partial\ SSR}{\partial\ b_3} \cdot \text{Learning Rate}$$

$$\text{New } b_3 = \text{Old } b_3 - \text{Step Size}$$

# Backpropagation

- Gradient descent with backpropagation is not guaranteed to find the global minimum of the error function, but only a local minimum; also, it has trouble crossing plateaus in the error function landscape.

- This issue, caused by the non-convexity of error functions in neural networks, was long thought to be a major drawback, but LeCun et al. (2015) argue that "recent theoretical and empirical results strongly suggest that local minima are not a serious issue in general."



LeCun, Yann, Yoshua Bengio, and Geoffrey Hinton. "Deep Learning." *Nature* 521, no. 7553 (2015): 436-444.

# References

- Activation Functions in Neural Networks by Towards Data Science

- Common Loss functions in machine learning by Towards Data Science

- Deep Learning Basics by Google Colab

- Deep Learning Tutorial for Beginners by Simplilearn

- Estimation of Neurons and Forward Propagation in Neural Net by Analytics Vidhya

- Learning Rate in Machine Learning by Deepchecks

- Neural Networks / Deep Learning by StatQuest

- Python Deep Learning Tutorial by Tutorials Point

- Softplus and softminus

# Until Next Class…

**What you need to do, what you will do**

# Upcoming Schedules

- **Feb 21 (Fri):** <span style="color:darkred">**Hands-on Assignment #1 Due**</span>

- **Feb 27 (Thu): Next Week**
  - Deep Learning: Basics (Remaining Part)
  - Short Session for Mid-term Preparation

- **Mar 6 (Thu):** <span style="color:darkred">**Mid-term Exam**</span>

- **Mar 13 (Thu)**
  - Feedback for Hands-on Assignment #1
  - Deep Learning: Advanced (Concept + Hands-on)
  - Hands-on Assignment #2 is out.

- **Mar 20 (Thu):** No Class (*Spring Recess*)

**Jaeung Sim**

Assistant Professor
School of Business, University of Connecticut
jaeung.sim@uconn.edu