



Assignment 1

Question 2

—

Abhinav Prakash

21CH10001

Course Code: CH620003

Problem Statement

Most viral outbreaks have a finite time period (short duration latency) during which the individual is infected but not yet infectious. Typically, this is modelled using the susceptible (S) – infected (I) – recovered (R) model. Of course there are many variations to the SIR model, for this problem, we stick to the basic version. The recovered fractions include the individual who are expired as well as healthy (after getting infected) too.

During the outbreak period, the total population (N) is assumed constant (through new births, migration, deaths, etc.), $\dot{N} = \dot{S} + \dot{I} + \dot{R} = 0$.

The dynamics of the disease transmission in this case is described as

$$\frac{dS}{dt} = -\frac{\beta}{N}SI; \frac{dI}{dt} = \frac{\beta}{N}SI - \gamma I; \frac{dR}{dt} = \gamma I. \quad (2)$$

Here, β represents the average contact rate of infection, signifying the transmission rate in the population. The average recovery rate is denoted by γ , and can be strongly correlated to the clinical observation for different virus variants. The reproduction number, $R_0 = \beta/\gamma$, which suggests when the transmission rate is more than the recovery, $R_0 > 1$, leads to an outbreak.

(a) Use the data available (from any online resource) for the COVID-19 spread dynamics in India and report the parameters R_0 and γ based on the above SIR model equation, for the first, second and possibly the third wave.

(b) Develop an optimised Neural network model to correlate the number of infected and recovered individuals with time from the initial period, starting from mid-March 2020 (you can ignore the few cases which were reported in Jan and Feb 2020). You may need to use the long short term memory (LSTM) time series network model or other recurrent network models.

Divide the total available dataset (available till today from March 2020) into training and testing data. The train dataset can be the first 'x' number of months (say 0-30 months, or the first wave as training and second wave as testing, or something similar) and remaining as test dataset. You should NOT randomly pick data points for the testing.

The mean squared error on the test data set should be more than 0.9. Note that adjusting the relative proportion of data available for training and testing, can drastically affect the model predictions over the test data.

(c) Find out the optimised number of nodes needed for the problem. You may consider only 1 hidden layer.

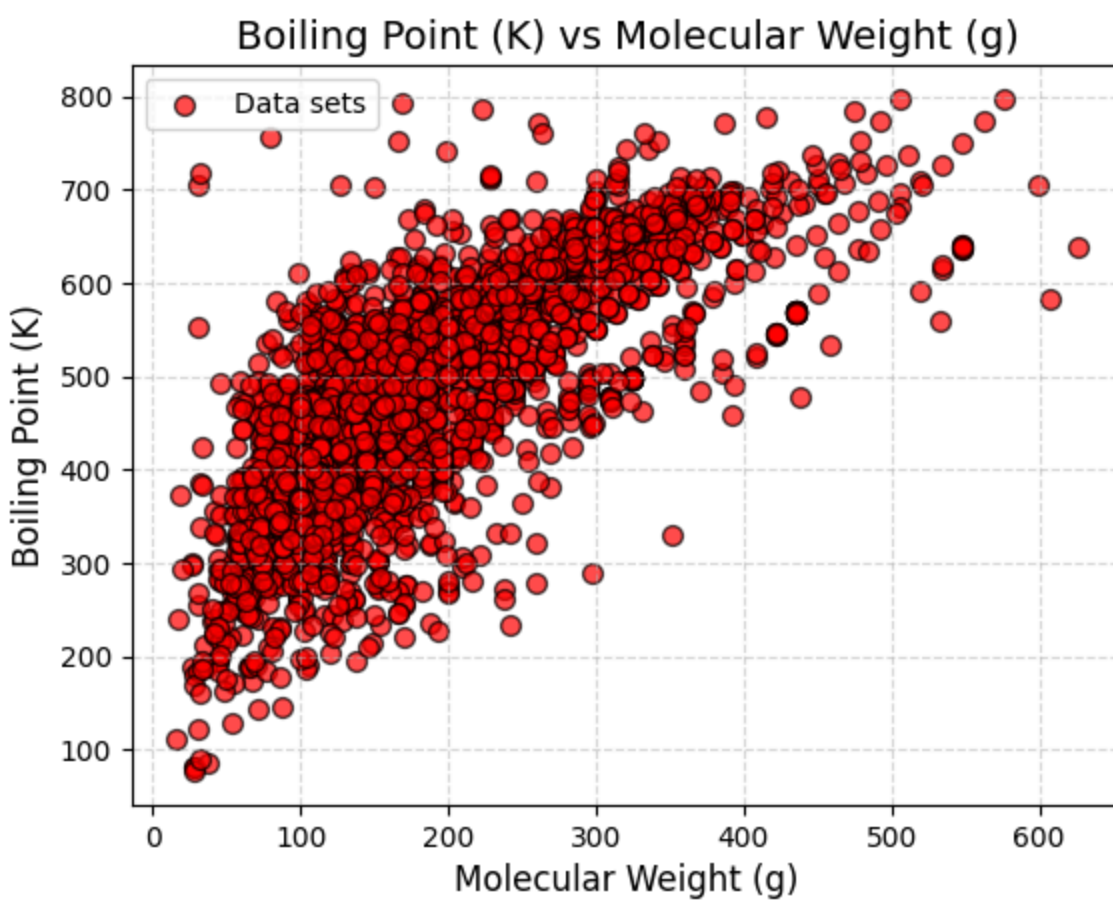
(d) Make a future forecast of the infection dynamics for the next 2 years, starting from today, using your trained Neural Network model.

(Hint: Please make some background reading about the LSTM model and use standard available codes for the same)

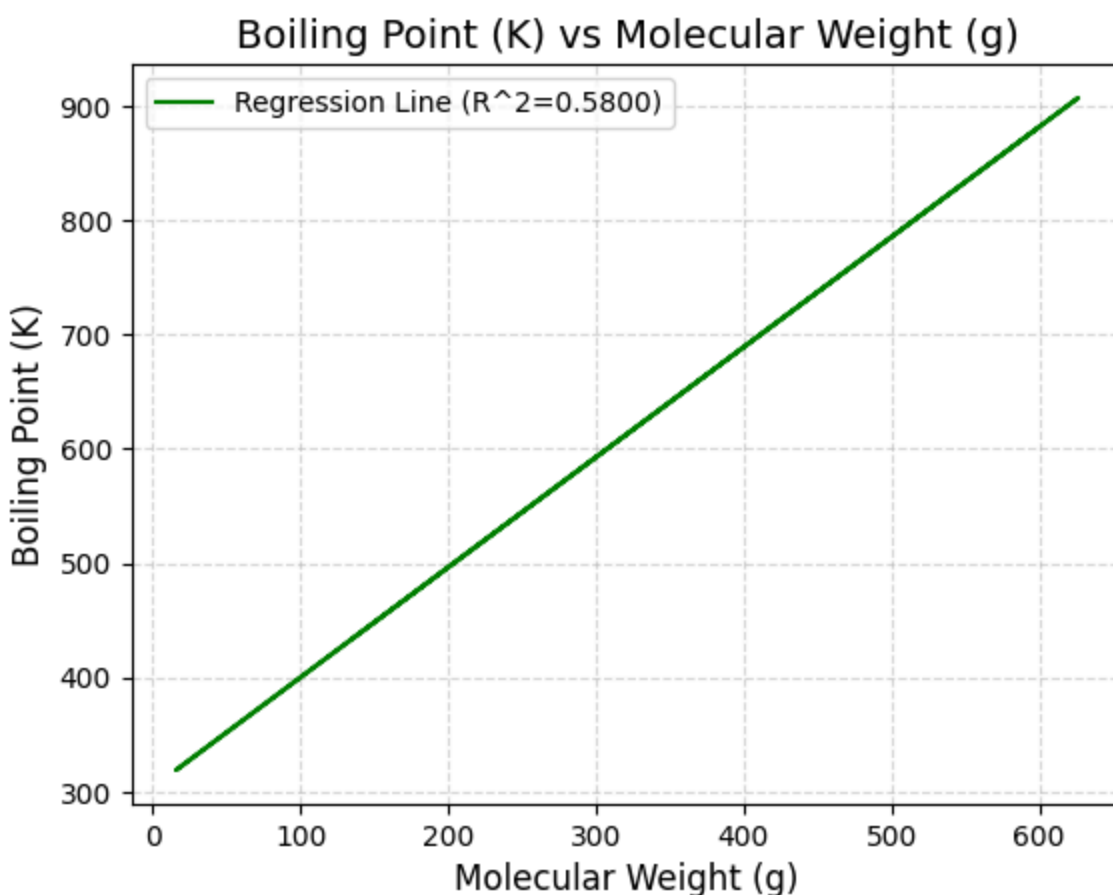
Google Colab Link

<https://colab.research.google.com/drive/1INoHybCifEAblRG7wsHaoQBWVCPSXuTH?usp=sharing>

Graph - BP vs Mw



Fit a straight line to corresponding scatter plot



The R2 value for the plotted line is notably low, measuring **below 0.6**. Consequently, a linear regression may not be deemed suitable for this dataset.

Further Data Processing

In my approach to data preprocessing, my objective was to extract three key columns: compound name, molecular weight, and acentric factor. To address chiral compounds identified by '+' and '-', I introduced a binary representation (0 for '+', 1 for '-'). Following this, I removed duplicate entries based on both the compound name and the chiral flag, ensuring that the dataset retained only unique values for subsequent analysis. The resultant dataset now features the following columns: Name (with corresponding float values), Molecular Weight, and Acentric Factor.

```
ds["name"] = new(ds["name"])
```

ds

	name	molweight	acentric factor	boiling point (K)
0	0.5	136.23704	0.3410	428.65
1	0.5	136.23704	0.2960	432.65
2	0.5	136.23704	0.3410	429.35
3	0.5	156.26820	0.6120	498.65
4	0.5	136.23704	0.2960	439.95
...
6026	0.5	18.01528	0.3449	373.15
6027	0.5	106.16740	0.3170	413.15
6028	0.5	410.84000	0.3100	633.15
6029	0.5	233.03480	0.2980	604.15
6030	0.5	598.84188	0.3430	704.15

6031 rows × 4 columns

Theta Values:

Theta0: 496.4208942392646

Theta1: 0.7933792601750705

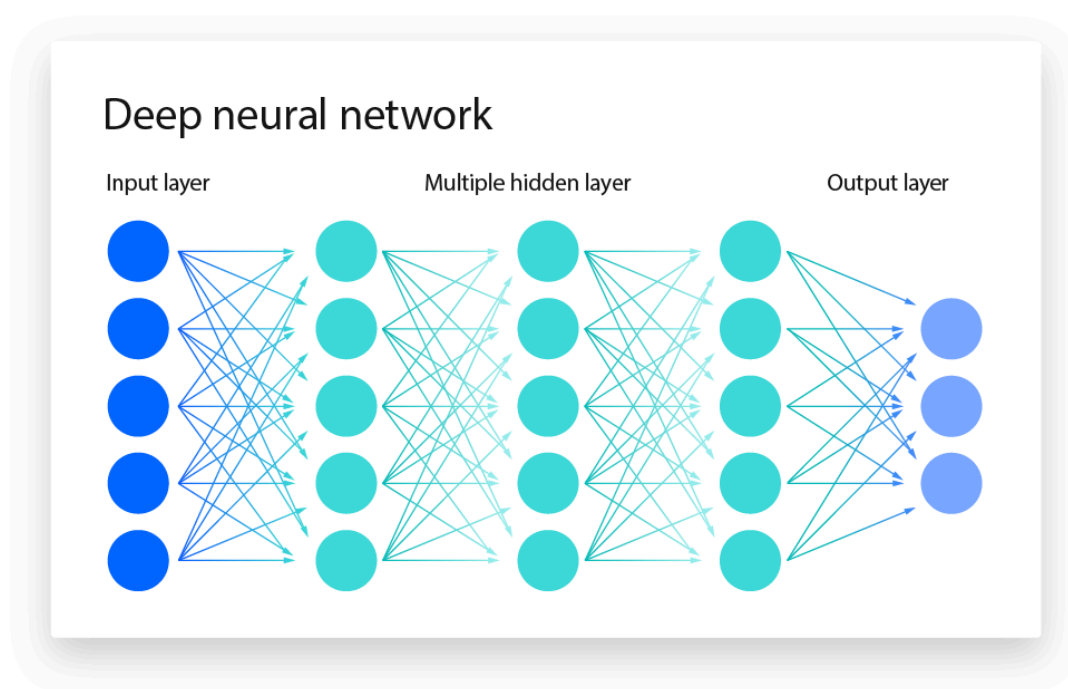
Theta2: 161.5616492950219

For the linear regression model, the following equation will be used:

Boiling point = Theta0 x Name + Theta 1 x molweight + Theta2 x acentric factor

ANN Theory:

Artificial Neural Networks (ANNs) are computational models inspired by the structure and functioning of the human brain. Comprising interconnected nodes or "neurons," organized into layers, ANNs excel in learning complex patterns and relationships within data. The network undergoes a training process where it adjusts its internal parameters to minimize the difference between predicted and actual outcomes. ANNs find widespread applications in various fields, including image and speech recognition, natural language processing, and predictive analytics, leveraging their ability to capture intricate non-linear dependencies in data for efficient decision-making and pattern recognition.



Input layer

The input layer in an Artificial Neural Network (ANN) represents the initial stage where raw data or features are introduced to the network. Each node in this layer corresponds to a specific feature, and the layer serves as the entry point for information into the neural network.

Hidden layer

A hidden layer in an ANN resides between the input and output layers, and its nodes collectively contribute to the network's ability to capture complex relationships within the data. These hidden layers enable the network to learn hierarchical representations, transforming input features into higher-level abstractions through weighted connections and activation functions.

Output layer

The output layer is the final layer in an ANN responsible for producing the network's output or predictions. The nodes in this layer generate the model's output based on the patterns and relationships learned during the training process. The number of nodes in the output layer depends on the nature of the task, such as binary classification, multi-class classification, or regression, with each node corresponding to a specific class or predicted value.

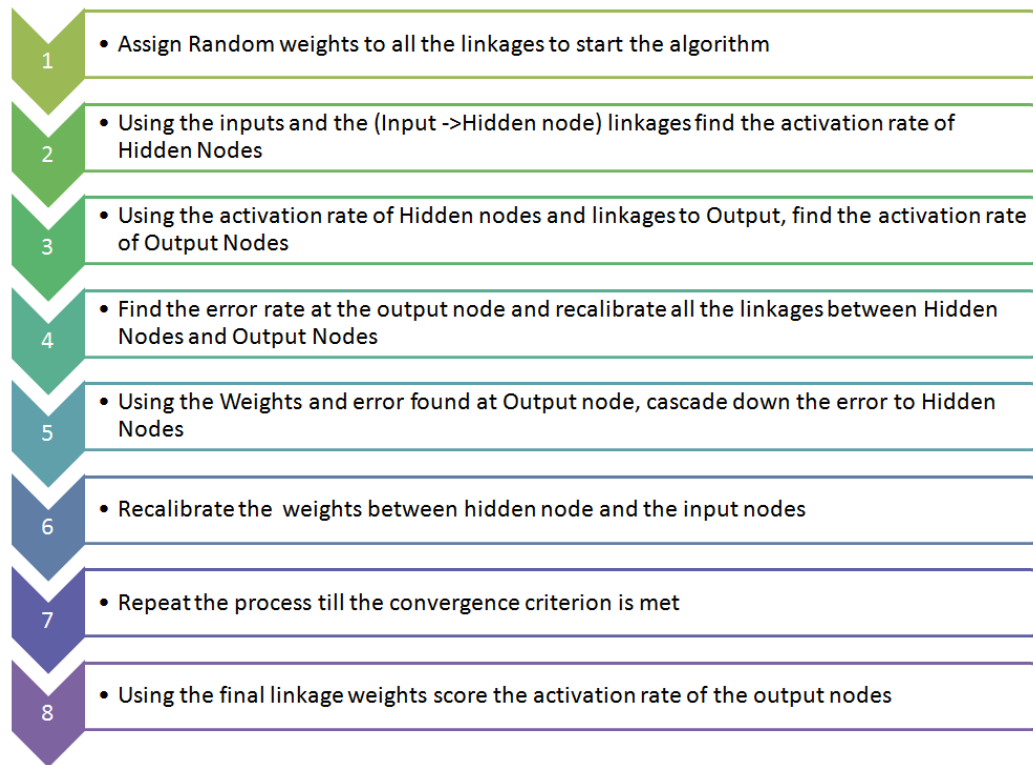
Activation Function

The **activation function** in an Artificial Neural Network (ANN) serves a crucial role in introducing non-linearity to the network, enabling it to learn complex relationships within data. It operates on the weighted sum of inputs and determines the output of a node or neuron. Common activation functions include the sigmoid, hyperbolic tangent (tanh), and rectified linear unit (ReLU). Sigmoid and tanh functions squash the output between specific ranges, making them suitable for binary classification and normalization tasks. ReLU, on the other hand, introduces non-linearity by allowing positive values to pass through and setting negative values to zero, contributing to the network's ability to handle intricate patterns and improve training efficiency. The choice of activation function depends on the nature of the problem and the desired characteristics of the network's learned representations.

ANN ARCHITECTURE:

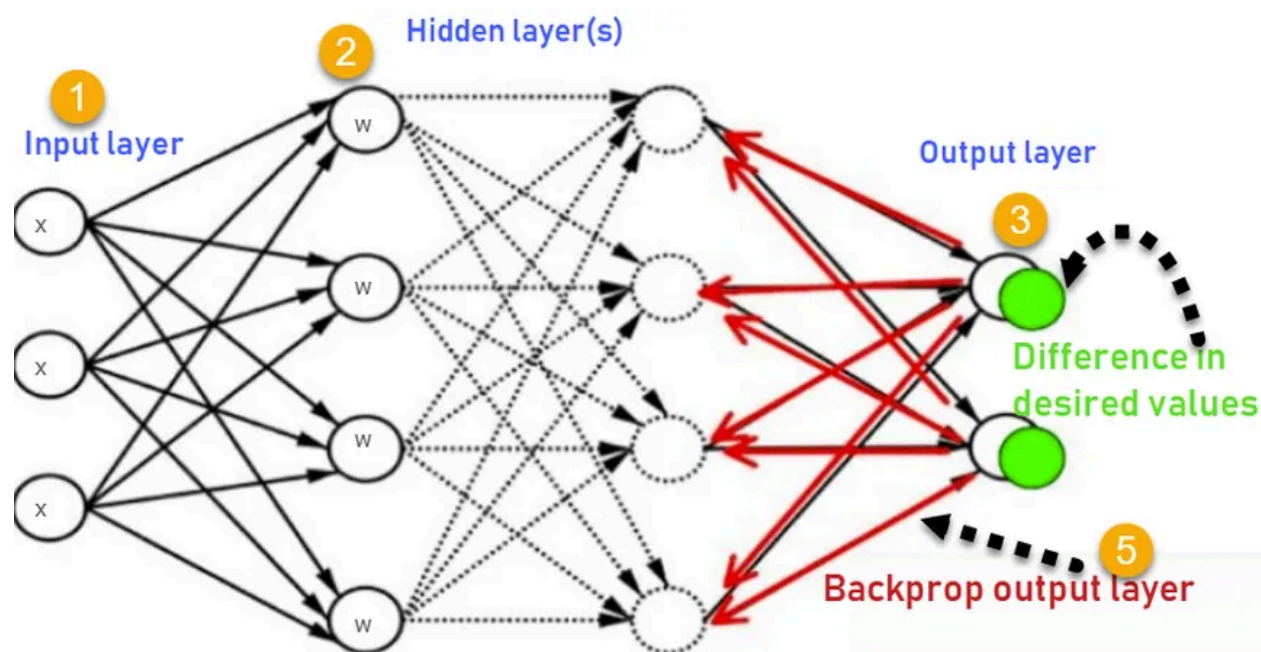
ANN Architecture:

The architecture of an Artificial Neural Network (ANN) defines its organizational structure, encompassing the arrangement and connectivity of its layers of neurons. A typical ANN architecture consists of three main components: the **input layer**, the **hidden layers**, and the **output layer**. The **input layer** receives the initial data, and each neuron in this layer represents a feature of the input. The **hidden layers**, positioned between the input and output layers, process and transform the input through weighted connections, applying activation functions to introduce non-linearity. Multiple hidden layers allow the network to learn hierarchical representations of data. The **output layer** produces the final result of the network's computation, providing the desired output based on the learned features. The number of neurons in each layer and the connectivity between them define the network's capacity to learn and generalize from the given data. The design and size of the architecture are crucial considerations, impacting the network's ability to comprehend complex patterns and perform effectively on diverse tasks.



Backward propagation:

Backpropagation, short for "backward propagation of errors," is a critical process in training Artificial Neural Networks (ANNs). Imagine you have a teacher guiding you on a task. You make a guess, and the teacher corrects your mistakes. Backpropagation works similarly for ANNs. When the network makes a prediction, we compare it to the actual result. If there's a difference (an error), we go backward through the network, adjusting the connection weights to reduce future errors. It's like learning from your mistakes and getting better at the task with each attempt. Backpropagation helps ANNs improve and fine-tune their predictions, making them more accurate over time.



ANN Build <Code>:

<Code>

```
import tensorflow as tf

# Define the model
ann = tf.keras.Sequential([
    tf.keras.layers.Dense(128, activation="relu",
input_shape=(2,)),
    tf.keras.layers.Dropout(0.2),
```

```

tf.keras.layers.BatchNormalization(),
tf.keras.layers.Dense(64, activation="relu"),
tf.keras.layers.Dropout(0.1),
tf.keras.layers.BatchNormalization(),
tf.keras.layers.Dense(units=1)
])

# model configure
optimizer = tf.keras.optimizers.Nadam(learning_rate=0.001)
ann.compile(loss="mse", optimizer=optimizer, metrics=["mae"])

# Train model
# Replace "features" and "labels" with your actual data
ann.fit(X_train, y_train, epochs=100, batch_size=32)

```

This piece of code does **two main things**: it gets a computer ready to learn from data and then teaches it. First, it takes information about molecular weight and acentric factor to predict boiling points. The data is split into parts for training (90%) and testing (10%). The computer is taught to understand this data through a neural network, which is like a virtual brain. This network has layers that help in learning patterns. The goal is to make accurate predictions about boiling points. The computer goes through this learning process 100 times, making adjustments each time to get better at making predictions. Overall, it's about making a smart system that can predict boiling points based on given information.

<Output>

```

170/170 [====] - 1s 5ms/step - loss: 2370.3655 - mae: 34.5597
Epoch 99/100
170/170 [====] - 1s 5ms/step - loss: 2370.3655 - mae: 34.5597
Epoch 100/100
170/170 [====] - 1s 5ms/step - loss: 2410.5371 - mae: 34.9046
<keras.src.callbacks.History at 0x7dd690b9bf10>

```

The Mean Absolute Error (MAE) of the improved neural network model has dropped to 34.9081, thanks to using dropout and batch normalization techniques. In the original model without these adjustments, the MAE was very high, around 104. It's important to mention that artificial neural networks (ANNs) might not work as well with smaller datasets, like the one here with only 6000 data points. For smaller datasets, other models might give better results, as ANNs could have trouble fitting the data too closely or making general predictions.

```
19/19 [=====] - 0s 2ms/step

|| Normalized Prediction : 462.573547, Actual: 419.5500 ||
|| Normalized Prediction : 488.632935, Actual: 476.6500 ||
|| Normalized Prediction : 548.997498, Actual: 607.1500 ||
|| Normalized Prediction : 517.453552, Actual: 489.7300 ||
|| Normalized Prediction : 473.697449, Actual: 463.1500 ||
-----
```

It's evident that our model is not capturing the complexity of the data well, indicating underfitting. To improve results, more data will be necessary.
