



Comparison of Cobb-Douglas Production Function with Deep Learning Models to Estimate the Productivity of Organised Manufacturing Sector - Indian Perspective

P. Preethi^{1*}, S. A. Jyothi rani², V.V. Haragopal³

¹Research Scholar, Department of Statistics, Osmania University, Hyderabad, India.

²Professor, Department of Statistics, Osmania University, Hyderabad, India.

³Professor (Retd.), Department of Statistics, Osmania University, Hyderabad, India.

Email: ²jyothi.rani@osmania.ac.in, ³haragopal_vajjha@yahoo.com

Corresponding Email: ^{1*}preethipatil24@gmail.com

Received: 17 May 2024

Accepted: 12 August 2024

Published: 28 September 2024

Abstract: *The production function defines the output of a firm, industry, or economy based on various combinations of inputs. This study focuses on estimating the production function of India's manufacturing sector and determining the relationship between its inputs and outputs. To predict the output, measured as Gross Value Added (GVA), the study applies the Cobb Douglas Production Function (CDPF) along with several deep learning techniques, including Feedforward, Recurrent, Long Short-Term Memory (LSTM), and Bidirectional LSTM networks. Model performance is evaluated using metrics like Mean Square Error (MSE), Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and Mean Absolute Percentage Error (MAPE). The results indicate that the Cobb Douglas Production Function outperforms the deep learning models in predicting GVA for the organized manufacturing sector in India.*

Keywords: *Gva, Mlr, Cdpf, Ffnn, Rnn, Lstm, Bidirectional Lstm.*

1. INTRODUCTION

The production function defines the resulting output of a company, an industry, or an entire economy across various input combinations. The performance of the manufacturing sector plays a pivotal role in anticipating the expansion of the industry sector, subsequently influencing India's overall gross domestic product (GDP) growth. Currently, 25% of the GDP share is contributed by the Industry sector, under that, the major share comes from the manufacturing sector only. The Annual Survey of Industries provides annual data on the organised manufacturing sector with different parameters like the number of factories established, fixed capital, total persons engaged, and so on.



In this research, we have chosen GVA (the dependent variable) to represent the sector's output, while eight other variables have been designated as independent variables. These include fixed capital (FC), total persons engaged (TPE), number of factories (N), material consumed (MC), interest paid (IP), total emoluments (TE), rent paid (RP), and fuel consumed (FUC).

For our analysis, we have employed a variety of models, namely the CDPF (Cobb-Douglas Production Function), FFNN (Feed-Forward Neural Network), RNN (Recurrent Neural Network), LSTM (Long Short Term Memory), and Bidirectional LSTM (Bidirectional LSTM).

2. RELATED WORK

Aliahmadi, Alireza, et al. (2016): “This paper uses both a multiple regression equation and an Artificial Neural Network (ANN) model to forecast total productivity based on labor, production, and industrial workshop productivity from 2006 to 2012. The two methods are compared using five criteria, and the results show that while the ANN model performs better than linear regression, the difference is not highly significant.” [1]

Kablay, Hassan, and Victor Gumbo: “This study applies Multiple Linear Regression (MLR) and a feed-forward Neural Network (NN) to predict the performance of 11 banks in Botswana, using Return on Assets (RoA) as the dependent variable and factors like management quality, credit risk, and liquidity as independent variables. Data from 2015-2019 financial reports is analyzed, and in the MLR model, the cost-to-income (C_I) ratio and loan loss provision to total loans (LLP_TL) ratio are identified as the most significant drivers of bank performance.” [2]

Santin, Daniel, Francisco J. Delgado, and Aurelia Valino (2004): “This paper demonstrates how artificial neural networks (ANN) serve as a valid semi-parametric alternative for fitting production functions and measuring technical efficiency. A Monte Carlo experiment is conducted using simulated production technology to compare ANN's efficiency results with Data Envelopment Analysis (DEA) and Corrected Ordinary Least Squares (COLS). Since ANN provides average production function estimates, the paper introduces a "thick frontier" strategy to convert these average estimations into a productive frontier.” [3]

Objectives of Study

- To calibrate the model utilizing the Cobb-Douglas approach
- To configure the model using Python code and various deep learning techniques (FFNN, RNN, LSTM, and Bidirectional LSTM).
- Ultimately, selecting the optimal model for estimation by evaluating its performance based on MAE, MAPE, MSE, and RMSE metrics.

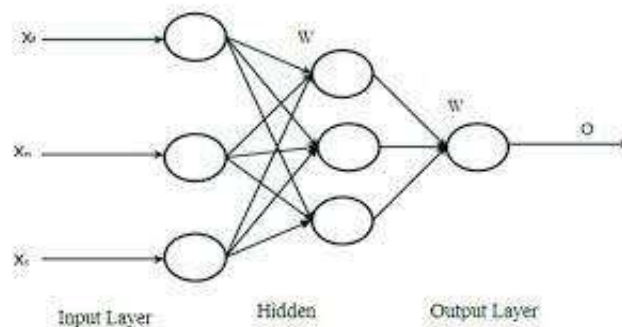
Data Source

“The study considers the Annual Survey of Industries data, published by the Ministry of Statistics and Programme Implementation (MOSPI) every year. This study examines the above-mentioned objectives at all Indian levels. The period of study is from 1981 to 2017” [4].
Software used: In analysing the data in this context we have used the Stata-17 and Python software package for finding the results.

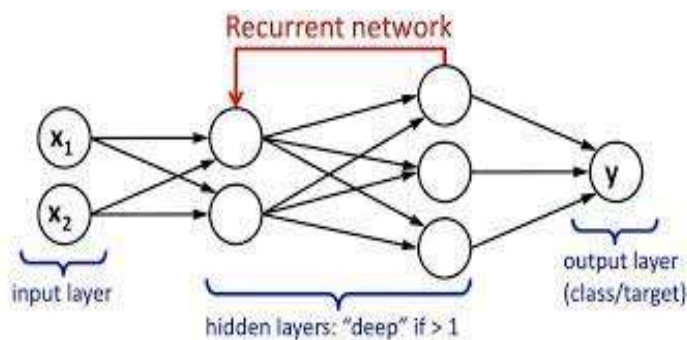
3. METHODOLOGY

Neural Network Architectures: In an Artificial neural network, the basic unit is neurons and there are many neurons that are interconnected to each other forming the network that is why it is called neural networks. There are different types of architectures in neural networks and a few of the most widely used are described below.

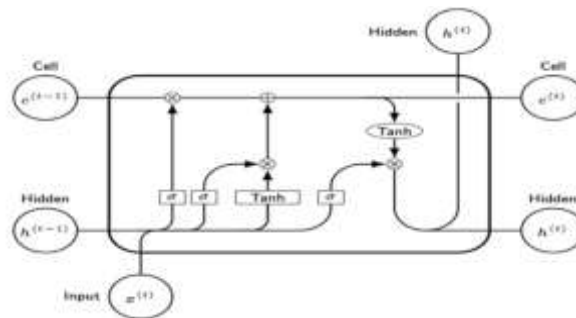
Feed Forward Neural Networks: Feed-forward neural networks direct information in a single direction, from input to output, without backward loops. They can be classified as single-layer or multi-layer based on the number of layers.



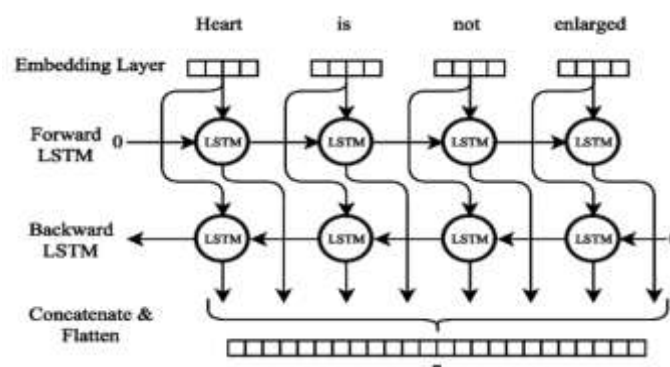
Recurrent Networks: Unlike feed-forward networks, recurrent neural networks (RNNs) allow information to flow in both forward and backward directions. The basic RNN topology is fully recurrent, where every neuron is connected to all others, including self-feedback loops.



Long Short-Term Memory (LSTM): LSTM, a special type of RNN, addresses short-term memory issues and learns long-term patterns by overcoming challenges like gradient vanishing, complex training, and processing long sequences. It is widely used in applications like music composition, speech recognition, and time series prediction.



Bidirectional Long-Short-Term Memory (BI-LSTM): It allows a neural network to process sequence information in both directions, from past to future and future to past, unlike regular LSTM, which flows in only one direction. This dual flow preserves both past and future information for better context understanding.



4. DISCUSSION AND RESULTS

4.1. Non-Linear Regression Model: Cobb-Douglas Model

“In its generalised form, the Cobb-Douglas function models for more than two factors may be written as” [5]

$$f(x) = A \prod_{i=1}^n x_i^{\lambda_i}, \quad x = (x_1, \dots, x_n).$$

Here, gross value added (GVA) is considered an output variable, and fixed capital (FC), total person engaged (TPE), number of factories (N), material consumed (MC), interest paid (IP), total emoluments (TE), rent paid (RP), and fuel consumed (FUC) are considered input variables.

$$GVA = A (FC)^{\alpha_1} (TPE)^{\alpha_2} (N)^{\alpha_3} (MC)^{\alpha_4} (IP)^{\alpha_5} (TE)^{\alpha_6} (RP)^{\alpha_7} (FUC)^{\alpha_8} \dots (1)$$

By applying log on both sides, we get

$$\text{Log GVA} = \text{Log A} + \alpha_1 \text{Log FC} + \alpha_2 \text{Log TPE} + \alpha_3 \text{Log N} + \alpha_4 \text{Log MC} + \alpha_5 \text{Log IP} + \alpha_6 \text{Log TE} + \alpha_7 \text{Log RP} + \alpha_8 \text{Log FUC}$$

$$Y = a + \alpha_1 X_1 + \alpha_2 X_2 + \alpha_3 X_3 + \alpha_4 X_4 + \alpha_5 X_5 + \alpha_6 X_6 + \alpha_7 X_7 + \alpha_8 X_8 \dots (2)$$



Where X1=Log FC X2= Log TPE X3= Log N X4= Log MC X5=Log IP X6=Log TE X7=Log RP X8=Log FUC

The above equation is in the form of a multiple regression equation. Therefore, the multiple regression model is fitted. Here the model is fitted by splitting 70% training and 30% testing data sets. Therefore the fitted model for 70% of the training dataset is shown below.

Table I. Fitting of Non-Linear Regression Model (Cobb-Douglas Production Function) for Training Data Set

Source	SS	df	MS	F	Pr > F	T	Pr > T	[95% Conf. Interval]
Model	22.1884788	8	2.7883861	2.10	0.044	0.097839	0.001756	
Residual	.021439632	16	.001339803			-.2714324	.7293111	
Total	22.1891183	24	.924171598			-.0285031	.5007642	

T	Coeff.	Std. Err.	T	Pr > T	[95% Conf. Interval]
X1	.6016413	.2441699	2.44	0.026	.0797846 1.123498
X2	.3474727	.1553099	2.10	0.044	-.0097839 .6951756
X3	-.0784394	.3045577	-0.26	0.801	-.6714324 .5145526
X4	.2361305	.1240320	1.89	0.077	-.0285031 .5007642
X5	-.2924338	.105129	-2.78	0.013	-.5122974 -.0695702
X6	.0443764	.2304866	0.20	0.849	-.4422333 .534886
X7	.1878819	.1186971	1.41	0.177	-.0839447 .4199086
X8	-.1724886	.1143718	-1.51	0.151	-.4099886 .0650114
_cons	-5.577304	3.842934	-1.45	0.166	-.13.72396 2.569321

From the above Table I it is observed that p value for t are less than α (0.10) for X1, X2, X4 and X5 therefore X1, X2, X4 and X5 have significant impact on output Y.

Therefore model equation formed from the relationship between dependent variable Y and 4 independent variables (X1, X2, X4 and X5) is obtained as

$$Y = -5.577304 + 0.6016413 X1 + 0.3474727 X2 + 0.2361305 X4 - 0.2924338 X5$$

4.2. Deep Learning Model: Ffnn, Rnn, Lstm Bi-Lstm

After studying the different Cobb-Douglas Model we now evaluate the deep learning models viz FFNN, RNN, LSTM and BI-LSTM.

Using Python software the different models were generated. First, different libraries were imported using import syntax as shown below. “NumPy is used to perform mathematical operations on arrays, pandas are used for working with data sets, seaborn making statistical graphics, and matplotlib is used for creating static, animated, and interactive visualizations, etc” [6].

```
import numpy as np
import pandas as pd
import seaborn as sns
from matplotlib import pyplot as plt
```

“By using `pd.read_excel` the data is loaded on the Python software” [7]

```
df = pd.read_excel('ASIneuralnetwork1.xlsx')
print(df.head())
```



After loading the data, it is divided into two distinct datasets: one for training and the other for testing. To determine the optimal split ratio for the training and testing datasets, various models are applied with different configurations. The findings reveal that MSE, MAE, RMSE and MAPE are minimized when 70% of the observations are allocated to the training set. To perform this data splitting, “the train test split function from the sklearn. Model selection library is utilized” [8]

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(x, y, test_size=0.3, random_state=42)
```

Hence, we allocate 70% of the data for training and 30% for testing across various models using,

- Activation function: Rectified Linear Unit
- Batch and Sequence Size: 1 and 12
- Epochs:100

“These models are implemented in Python, employing the Keras library, renowned for its ease of use and effectiveness in developing and assessing deep learning models” [9]

```
from keras.models import Sequential
from keras.layers import Dense, Activation, Dropout, Input, InputLayer, SimpleRNN, LSTM, Bidirectional
```

The model is executed with different input and dense layers then it is fitted by using model. Fit command in Python.

```
history = model.fit(X_train, y_train, verbose=1, epochs=100, validation_data=(X_test, y_test))
```

“Utilising MSE, MAE, RMSE and MAPE, various approaches are applied to the data” [10].

$$\text{MSE} = 1/n \sum_{t=1}^n (Y_t - F_t)^2$$

$$\text{MAE} = 1/n \sum_{t=1}^n |Y_t - F_t|$$

$$\text{MAPE} = \frac{1}{n} \sum_{t=1}^n \left| \frac{Y_t - F_t}{Y_t} \right| \times 100$$

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{t=1}^n (Y_t - F_t)^2}$$

Where n is the number of years, F_t is the estimated or predicted value, Y_t is the actual value. Table II presents a comprehensive overview of various model specifications, including their names, training dataset percentages, and key performance metrics such as MSE, MAE, RMSE, and MAPE. A clear pattern emerges from the data: the Cobb-Douglas Production Function consistently outperforms other methods in terms of MSE, MAE, RMSE, and MAPE. This suggests that implementing the Cobb-Douglas Production Function for future forecasting holds promise, as it consistently exhibits lower error rates compared to alternative methods.

Table II. Test Accuracy

Model Name	Specification	Training dataset	MSE	MAE	RMSE	MAPE
Cobb-Douglas Production Function	-	70%	0.0005	0.0423	0.0470	0.2531
FFNN	DL(64) DL(32) DL(1)	70%	0.0044	0.1419	0.1765	0.8273
RNN	Simple RNN(64) DL(1)	70%	0.0018	0.1295	0.1706	0.7517
LSTM	LSTM(50) LSTM(50) DL(1)	70%	0.0076	0.2000	0.2434	1.1876
Bidirectional LSTM	Bidirectional (LSTM(50) DL(1)	70%	0.0077	0.1650	0.2094	0.9564

We now discuss the validation process of all the methods. The loss values during training and validation offer a more profound understanding of how the learning performance evolves across epochs. They also assist in identifying issues with the learning process, which can result in either an underfit or overfit model. The Figure 1 depicts that the value of epochs equal to 100 is a good fit and there is no problem with the learning rate.

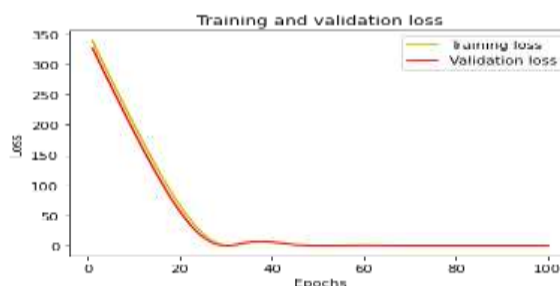


Figure 1. Training and Validation Loss Curve

The graphs are also plotted for the same dataset for all the models (CDPF, FNN, RNN, LSTM, BI-LSTM) for training and testing datasets using plt command for plotting figure, line, etc

```
def plot_result(x_train, y_test, prediction_train, prediction_test):
    actual = np.append(y_train, y_test)
    predictions = np.append(prediction_train, prediction_test)
    rows = len(actual)
    plt.figure(figsize=(15, 10))
    plt.plot(y_train, actual)
    plt.plot(y_test, prediction_test)
    plt.plot(y_test, actual)
    plt.legend(['Actual', 'Predictions'])
    plt.xlabel('Years')
    plt.ylabel('GDP')
    plt.title('Actual and Predicted Values. The red line separates the training and testing datasets.')
    x_ticks = [2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, 2023, 2024]
    x_labels = range(0, rows)
    plt.xticks(x_ticks)
    plt.xticks(rotation=45, label='year')
```

The Figure 2 depicts the graph of Cobb-Douglas Model where the values before the redline indicate the training values and values after the redline indicate the testing values. Both the actual values and predicted values are fitted and the variation can be observed. Similarly in Figures 3, 4, 5, and 6 for all other selected models.



Figure 2. Cobb-Douglas Model

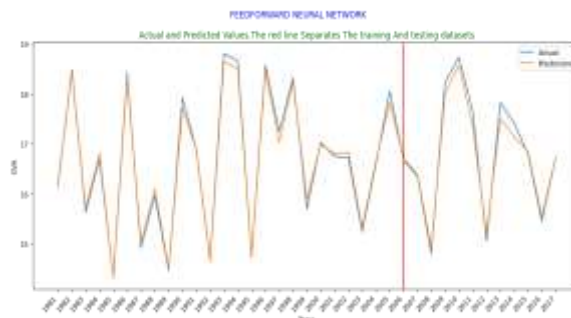


Figure 3. Feed Forward Neural Network

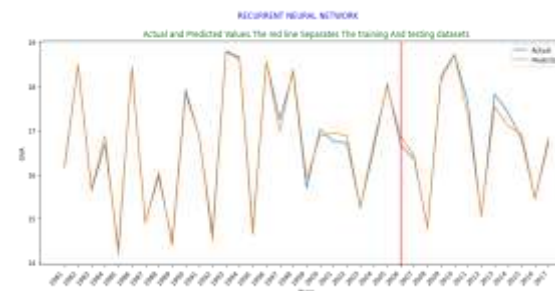


Figure 4. Recurrent Neural Network



Figure 5. LSTM



Figure 6. Bidirectional LSTM

Out of all the plotted figures almost all models are better fitting the data but the Figure 2 clearly shows Cobb-Douglas Model best fit for the selected data.

5. CONCLUSION

The test of accuracy (Table II) and the above-plotted graph (Figure 2) show that the Cobb-Douglas Model outperforms better than other models viz FFNN, RNN, LSTM, and Bidirectional LSTM. This clearly indicates that the Cobb-Douglas model performs well for the organised manufacturing sector. All other deep learning methods seem to perform very badly. Therefore, we used the Cobb-Douglas Model to fit considering 100% dataset to be a training dataset without considering the testing aspect, and the results are shown below in Table III.

Table III. Fitting of Cobb-Douglas Model

Source	SS	df	MS	Number of obs	F	Pr > F
Model	88.2880028	8	11.0250004	36	10.0000	<.0001
Residual	1.0727000	28	0.3831071			
Total	89.3607028	36				

V	Coefficient	Std. Err.	t	Pr > t	[95% Conf. Interval]
X1	0.1010131	0.0000000	1.07	0.288	-.0000000 .0000000
X2	0.2925455	0.0000000	2.92	0.006	0.0000000 .0000000
X3	-0.341566	0.0000000	-3.42	0.001	-0.0000000 -0.0000000
X4	0.2080857	0.0000000	2.12	0.042	0.0000000 0.0000000
X5	-0.3170303	0.0000000	-3.17	0.002	-0.0000000 -0.0000000
X7	0.2877669	0.0000000	2.88	0.007	0.0000000 0.0000000
X8	0.1922274	0.0000000	1.92	0.064	0.0000000 0.0000000
_cons	0.1010131	0.0000000	1.07	0.288	-0.0000000 0.0000000

From the above Table III, it is observed that p-value for t are less than α (0.10) for X2, X3, X4, X5, X7, and X8, therefore, X2, X3, X4, X5, X7, and X8 have a significant impact on output Y.

Therefore, the Cobb-Douglas equation model for the whole dataset can be written as
 $Y = 0.1010131 + 0.2925455 X2 - 0.341566 X3 + 0.2080857 X4 - 0.3170303 X5 + 0.2877669 X7 + 0.1922274 X8$



The above model can be rewritten in terms of Non-Linear Cobb-Douglas production as
 $GVA = 1.26 (TPE)^{1.96} (N)^{0.45} (MC)^{1.61} (IP)^{0.48} (RP)^{1.93} (FUC)^{1.55}$

The preceding model suggests an enhanced scale of returns. Furthermore, a 1% increase in total persons engaged (TPE), number of factories (N), material consumed (MC), interest paid (IP), rent paid (RP), and fuel consumed (FUC) corresponds to a respective increase of 1.96%, 0.45%, 1.61%, 0.48%, 1.93%, and 1.55% in gross value added (GVA).

Additionally, the findings reveal that both the Cobb-Douglas and Deep Learning Models effectively predicted the production function (Gross Value Added) for the specified period. Nevertheless, scrutiny of accuracy (see Table II) and the accompanying graph (refer to Figure 2) demonstrates that the Cobb-Douglas Model outperformed all other deep learning models.

6. REFERENCE

1. Aliahmadi, Alireza, et al. "Comparing linear regression and artificial neural networks to forecast total productivity growth in Iran." *International Journal of Information, Business and Management* 8.1 (2016): 93.
2. Kablay, Hassan, and Victor Gumbo. "Comparison of Multiple Linear Regression and Neural Network Models in Bank Performance Prediction in Botswana." (2021).
3. Santin, Daniel, Francisco J. Delgado, and Aurelia Valino. "The measurement of technical efficiency: a neural network approach." *Applied Economics* 36.6 (2004): 627-635.
4. MoSPI: Ministry of Statistics and Programme Implementation, 2021.
5. Peter Pesala, Indian manufacturing industry in the Era of globalization: a Cobb-Douglas production function analysis, *Indian Journal of Economics and Development* 5(1) (2017), 1-11.
6. Popular Python Libraries - NumPy, Pandas, Seaborn, Sklearn." *AlmaBetter*, 22 June 2023.
7. Pandas. Read_Excel — Pandas 2.2.1 Documentation. pandas.pydata.org/pandas-docs/stable/reference/api/pandas.read_excel.html.
8. Singh, Vikash Kumar, et al. "Impact of Train/Test Sample Regimen on Performance Estimate Stability of Machine Learning in Cardiovascular Imaging." *Scientific Reports*, vol. 11, no. 1, July 2021.
9. Rathod, Dushyantsinh B., and Vijaykumar Gadhavi. "A Review on Deep Learning Using Python Libraries and Packages." by *International Journal of Creative Research Thoughts (IJCRT)*, vol. 6, no. 2, May 2018, pp. 83–85.
10. Trevisan, Vinicius. "Comparing Robustness of MAE, MSE and RMSE - Towards Data Science." *Medium*, 25 Mar. 2022, towardsdatascience.com/comparing-robustness-of-mae-mse-and-rmse-6d69da870828.