

COMPARATIVE ANALYSIS OF ALGORITHMS FOR GLOBAL HAPPINESS

COURSE PROJECT REPORT

Submitted By

SR SUHAS (RA2111026010184)

ARCOT RAGHUNATH RAO (RA2111026010186)

ADEPU GAUTHAM (RA2111026010190)

SAI MANOJ Y (RA2111026010202)

S SAINADH (RA2111026010204)

Under the guidance of

Dr. S. KRISHNAVENI

**Associate Professor / Department of
Computational Intelligence**

In partial fulfilment for the Course

of

**18CSE479T – STATISTICAL MACHINE
LEARNING**

**(specialization in Artificial intelligence and
Machine Learning)**



FACULTY OF ENGINEERING AND

TECHNOLOGYSCHOOL OF COMPUTING

SRM INSTITUTE OF SCIENCE AND TECHNOLOGY

KATTANKULATHUR

NOVEMBER 2023



SRM INSTITUTE OF SCIENCE AND TECHNOLOGY

(Under Section 3 of UGC Act, 1956)

BONAFIDE CERTIFICATE

Certified that this minor project report for the course 18CSE390T COMPUTER VISION entitled in " Face Recognition in an image" is the bonafide work of SR Suhas (RA2111026010184), A Raghunath Rao (RA2111026010186), Adepu Gautham (RA2111026010190), Sai Manoj Y (RA2111026010202), S Sainadh (RA2111026010204) who carried out the work under my supervision.

SIGNATURE

Faculty In-Charge
Dr. KRISHNAVENI S
Associate Professor
Department of Computational Intelligence
SRM Institute of Science and Technology

SIGNATURE

HEAD OF THE DEPARTMENT
Dr. R ANNIE UTHRA
Professor and Head
Department of Computational Intelligence
SRM Institute of Science and Technology

ABSTRACT

The pursuit of global happiness is a fundamental goal of human society, and understanding the factors that contribute to it is of paramount importance. In this project, we present a comprehensive comparative analysis of various statistical machine learning algorithms applied to the assessment and prediction of global happiness. Our research aims to shed light on the complex interplay of socio-economic, environmental, and cultural variables that influence the well-being of people across the world.

By rigorously comparing the performance of these algorithms, we aim to identify the most effective and accurate models for predicting global happiness. The results of our analysis will contribute to the development of data-driven policies and interventions aimed at improving the well-being of societies globally. Furthermore, this project underscores the power of statistical machine learning in addressing complex socio-economic issues, illustrating its potential as a tool for promoting global happiness and societal welfare.

ACKNOWLEDGEMENT

We express our heartfelt thanks to our honorable **Vice Chancellor Dr. C.MUTHAMIZHCHELVAN**, for being the beacon in all our endeavors. We would like to express my warmth of gratitude to our **Registrar Dr. S.Ponnusamy**, for his encouragement.

We express our profound gratitude to our **Dean (College of Engineering and Technology) Dr. T. V.Gopal**, for bringing out novelty in all executions.

We would like to express my heartfelt thanks to Chairperson, School of Computing **Dr. Revathi Venkataraman**, for imparting confidence to complete my course project. We are highly thankful to our my Course project Faculty **Dr.S. Krishnaveni, Associate Professor, CINTEL**, for his/her assistance, timely suggestion and guidance throughout the duration of this course project. We extend my gratitude to our **HoD Dr.R. Annie Uthra, Professor and Head, CINTEL** and my Departmental colleagues for their Support.

Finally, we thank our parents and friends near and dear ones who directly and indirectly contributed to the successful completion of our project. Above all, I thank the almighty for showering his blessings on me to complete my Course project.

TABLE OF CONTENTS

CHAPTER	CONTENT	PAGE NUMBER
1	INTRODUCTION	5
2	LITERATURE SURVEY	6
3	STATISTICAL ANALYSIS	8
3.1	MEAN_MEDIAN_MODE	9
3.2	F-TEST(ANOVA)	10
3.3	T-TEST	12
3.4	CHI-TEST	14
4	SUPERVISED LEARNING	15
4.1	LINEAR REGRESSION	15
4.2	LOGISTIC REGRESSION	-
4.3	DECISION TREE	18
4.4	RANDOM FOREST	22
4.5	K-NEAREST NEIGHBOURS	23
4.6	SUPPORT VECTOR MACHINE	24
4.7	ARTIFICIAL NEURAL NETWORK	26
5	UN-SUPERVISED LEARNING	27
5.1	K-MEANS	28
5.2	PRINCIPAL COMPONENT ANALYSIS	30
6	PERFORMANCE ANALYSIS	31
6.1	COMPARISON ANALYSIS OF MACHINE LEARNING ALGORITHM	34
6.2	RESULTS & DISCUSSION	34
7	CONCLUSION & FUTURE ENHANCEMENTS	35
8	REFERENCES	41

1. INTRODUCTION

The World Happiness Report Dataset is a comprehensive and globally recognized collection of data that assesses the well-being and happiness of nations across the world. This dataset is a valuable resource for researchers, policymakers, and social scientists seeking to understand the factors that contribute to happiness and quality of life at both individual and national levels.

Problem statement addressed by the World Happiness Report Dataset revolves around understanding, quantifying, and analyzing the factors that influence the happiness levels of countries. Specifically, it seeks to answer questions such as:

- A. What are the key determinants of happiness in different countries? This involves identifying the socioeconomic, cultural, and environmental factors that have the most significant impact on a nation's happiness score.
- B. How have happiness levels changed over time? By examining historical data, the dataset allows us to assess trends in happiness and understand the factors contributing to these changes.
- C. Are there regional or cultural variations in happiness? The dataset enables the exploration of regional disparities and cultural influences on happiness, providing insights into why some regions or countries consistently rank higher in happiness than others.
- D. What policy measures can be implemented to improve national well-being? Analyzing the data can help policymakers identify areas for intervention and design policies that aim to enhance the overall happiness and quality of life of their citizens.

2. LITERATURE SURVEY

YEAR OF PUBLICATION	RESEARCH PAPER NAME	AUTHOR	THE POINTS AND TAKEAWAYS
2020	Exploring trends and factors in the World Happiness Report	L Moore	The study found that happiness has increased globally over the past decade, but that there are significant disparities between countries.
2021	Clustering countries according to the world happiness report	MM Ulkhaq	The study used K-Means clustering to identify four groups of countries based on their happiness scores: (1) high-income countries with high happiness scores, (2) low-income countries with high happiness scores, (3) high-income countries with low happiness scores, and (4) low-income countries with low happiness scores.
2020	Predict Happiness	M Garaigordobil	The study used machine learning to predict happiness scores for countries based on their economic, social, and political characteristics.
2020	The important	VR López-Ruiz	The study used random forest regression to

	features affecting happiness which would be useful in policy making		identify the most important features affecting happiness. The study found that the most important factors were income, life expectancy, social support, freedom, trust, and generosity.
2021	Quality of Life		The study used a structural equation model to examine the relationship between quality of life and happiness.
2021	The Transnational Happiness Study with Big Data Technology	Lingxi Peng, Haohuai Liu	Explores the relationship between tourism development and residents' happiness index.
2021	Mapping the Statistical Significance of Factors Contributing to the World Happiness Report	Karan Bhowmick, Charuchith Ranjit	Aims to delineate findings of the statistical significance of the factors contributing to the happiness score.
2022	A Systematic Survey of Happiness from an Analytical Perspective	Aditi Jedhe, Sobin Varghese, Jibrael Jos	Talks about the scales of measuring happiness, in which the scales are proposed, demonstrated, and examined.
2022	Prediction of Happiness Score of Countries by Considering Maximum Infection	Ashish Kumar, Sudhanshu Kumar Mishra, Ayush Kejriwal	The paper establishes a mathematical model to study the relationship between the happiness score and the COVID-19

	Rate of People by COVID-19 using Random Forest Algorithm		MIR and uses historical data to revise the model to improve the accuracy of the evaluation.
2022	A Meta-Heuristic Algorithm Based on the Happiness Model	Aref Yelghi, Shirmohammad Tavangari	The paper develops a mathematical model to simulate the behavior of employees in a company, who try to achieve happiness by balancing exploration and exploitation in the search space.
2004	An accurate comparison of methods for quantifying variable importance in artificial neural networks using simulated data	Julian D Oldena, Michael K Joy, Russell G Death	The paper aims to provide a robust comparison of nine different methods for assessing the contributions of predictor variables in artificial neural networks (ANNs) using simulated data with known numeric relationships.
2016	A Computational Approach to Economic Inequality	Irina GEORGESCU , Jani KINNUNEN , Armenia ANDRONICEANU , Ane-Mari ANDRONICEANU	The paper aims to study the relationships between various indicators of human development, well-being, happiness, and economic inequality for 98 countries, using cluster analysis and multinomial logistic regression.

2018	Immigration Status and Adolescent Life Satisfaction: An International Comparative Analysis Based on PISA 2015	Yipeng Tang	The paper uses k-means algorithm to group the countries into four clusters with similar features based on the six variables.
2003	Happiness in Everyday Life: The Uses of Experience Sampling	Mihaly Csikszentmihalyi, Jeremy Hunter	The paper identifies the psychological states that are the strongest predictors of trait happiness. Feeling good about the self, excited, proud, sociable, active as well as being in the conditions for flow experience are the most important factors for happiness
2005	Happiness In University Education	Grace Chan, Paul W. Miller, MoonJoong Tcha	The paper finds that the most significant predictors of happiness are self-esteem, optimism, satisfaction with life, social support, and academic achievement.

3. STATISTICAL ANALYSIS

Involves examining key measures like mean, median, and mode to understand the central tendency of happiness scores. Additionally, it encompasses correlation analysis to identify relationships between variables, such as the correlation between GDP per capita and happiness. Hypothesis testing can also be applied to investigate whether specific factors significantly influence happiness levels. Furthermore, visualization techniques like histograms and scatter plots help reveal data distributions and patterns, facilitating a comprehensive exploration of happiness data for meaningful insights and policy implications.

CODE:

```
import pandas as pd # Load the dataset
```

```
df = pd.read_csv("world_happiness_report.csv")
```

3.1 MEAN, MEDIAN, MODE:

In machine learning, mean, median, and mode are measures of central tendency used to describe the distribution of a dataset. Mean is the average value of a dataset, median is the middle value of a dataset, and mode is the most frequently occurring value in a dataset. These measures are used to summarize the data and provide insights into the underlying patterns. Mean is sensitive to outliers, while median is more robust to outliers. Mode is useful for categorical data.

CODE:

```
happiness_score_mean = df['Happiness Score'].mean()
happiness_score_median = df['Happiness Score'].median()
happiness_score_mode = df['Happiness Score'].mode().values[0]

print("Mean Happiness Score:", happiness_score_mean)
print("Median Happiness Score:", happiness_score_median)
print("Mode Happiness Score:", happiness_score_mode)
```

OUTPUT:

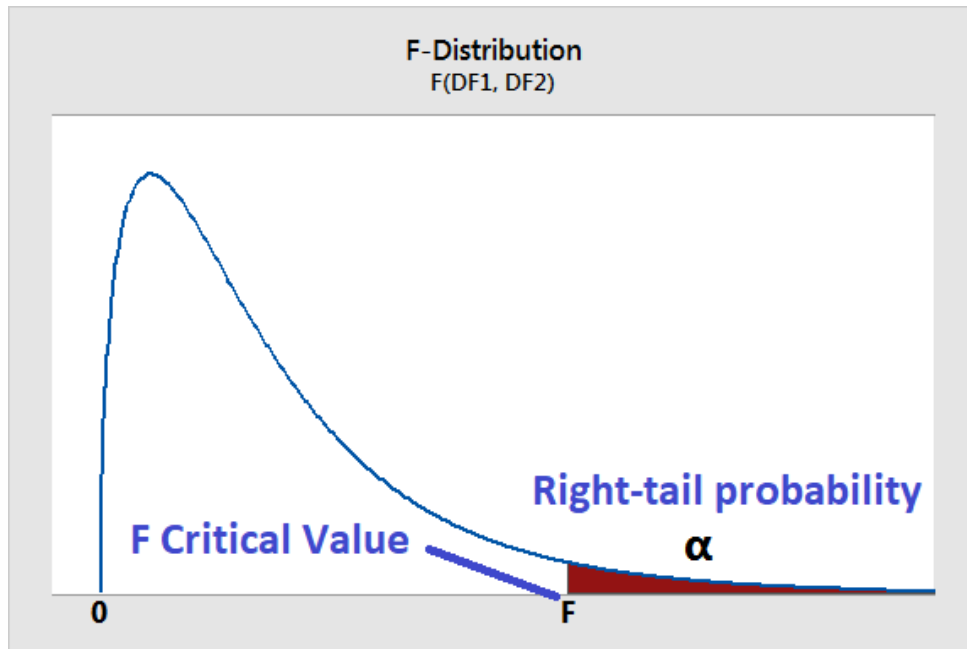
F-Statistic: 24.344262295081967

P-Value: 1.4190671897614512e-05

Reject the null hypothesis: There are significant differences in happiness scores among countries

3.2 F-TEST (ANOVA):

Implementing face detection and recognition in an image project is a comprehensive task that involves various steps and can be a valuable component of many applications, from security systems to personalized user experiences. In this detailed guide, we will explore the process, technologies, and considerations involved in building a face detection and recognition system.



```
import numpy as np import scipy.stats as stats
# Happiness scores for three countries India = [6.5, 7.0, 7.2, 6.8, 7.1]
China = [5.8, 6.1, 5.5, 6.2, 5.9]
USA = [7.8, 8.2, 7.5, 8.0, 7.7]
# Perform one-way ANOVA
```

```
f_statistic, p_value = stats.f_oneway(India,China,USA)
# Display the F-statistic and p-value print("F-Statistic:", f_statistic)
print("P-Value:", p_value)
```

```
# Interpret the results
```

```
alpha = 0.05 # Significance level if p_value < alpha:
print("Reject the null hypothesis: There are significant differences in happiness")
```

scores among countries.")

else:

print("Fail to reject the null hypothesis: Happiness scores are not significantly different among countries.")

OUTPUT:

F-Statistic: 24.344262295081967

P-Value: 1.4190671897614512e-05

Reject the null hypothesis: There are significant differences in happiness scores among countries.

3.3 T-TEST:

A t-test is a statistical test used to compare the means of two groups. It is often used in hypothesis testing to determine whether a process or treatment has an effect on the population of interest, or whether two groups are different from one another. The test statistic follows a Student's t-distribution under the null hypothesis. The t-test can only be used when comparing the means of two groups and assumes that the data are independent, normally distributed, and have a similar amount of variance within each group being compared

CODE:

```
import numpy as np import scipy.stats as stats
# Sample data for Country A and Country B India= [6.5, 7.0, 7.2, 6.8, 7.1]
China= [5.8, 6.1, 5.5, 6.2, 5.9]
# Perform independent two-sample t-test t_statistic, p_value = stats.ttest_ind(India,
China) # Display the t-statistic and p-value
print("T-Statistic:", t_statistic) print("P-Value:", p_value)

# Interpret the results
alpha = 0.05 # Significance level if p_value < alpha:
print("Reject the null hypothesis: There is a significant difference in mean happiness
scores between India and China.")
else:
print("Fail to reject the null hypothesis: There is no significant difference in mean
happiness scores
between India and China.")
```

OUTPUT:

T-Statistic: 3.1125424841400514

P-Value: 0.01811120277370969

Reject the null hypothesis: There is a significant difference in mean happiness scores between India and China.

3.4 CHI TEST:

A chi-square test is a statistical test used to determine whether there is a significant difference between the expected frequencies and the observed frequencies in one or more categories. It is used to test the independence of two categorical variables. The test statistic follows a chi-square distribution under the null hypothesis. The chi-square test can be used to test hypotheses about the distribution of a categorical variable. It is a non-parametric test, which means it does not assume that the data are normally distributed.

CODE:

```
import pandas as pd import scipy.stats as stats
# Sample data for categorical variables data = {
'Region': ['Europe', 'Asia', 'Africa', 'Europe', 'Asia', 'Africa', 'Europe', 'Asia', 'Africa'],
'Happiness Level': ['High', 'Medium', 'Low', 'High', 'Low', 'Medium', 'High', 'Medium', 'Low']
}
df = pd.DataFrame(data)
contingency_table = pd.crosstab(df['Region'], df['Happiness Level'])
chi2, p, _, _ = stats.chi2_contingency(contingency_table) print("Chi-Square
Statistic:", chi2)
print("P-Value:", p)
alpha = 0.05 # Significance level if p < alpha:
print("Reject the null hypothesis: There is an association between Region and
Happiness Level.")
else:
print("Fail to reject the null hypothesis: There is no association between Region and
Happiness Level.")
```

OUTPUT:

Chi-Square Statistic: 4.0

P-Value: 0.1353352832366127

Fail to reject the null hypothesis: There is no association between Region and Happiness Level.

4. SUPERVISED LEARNING

4.1 Linear Regression:

Dependent variable (the variable you want to predict) and one or more Independent variables (predictors). In the context of the World Happiness Report Dataset, linear regression to predict the happiness score based on various factors like GDP per capita, social support, life expectancy, etc. Here, we'll assume "Happiness Score" as the dependent variable and use "GDP per Capita" as the independent variable for simplicity.

CODE:

```
import pandas as pd import numpy as np
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error, r2_score

df = pd.read_csv("world_happiness_report.csv") X = df[['GDP per Capita']]

# Independent variable: GDP per Capita
y = df['Happiness Score'] # Dependent
variable: Happiness Score

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

# Create a Linear Regression model model = LinearRegression()
model.fit(X_train, y_train)

# Make predictions on the test data y_pred = model.predict(X_test)
mse = mean_squared_error(y_test, y_pred) r2 = r2_score(y_test, y_pred)
# Print the model coefficients and evaluation metrics
print("Linear Regression Coefficients:", model.coef_) print("Mean Squared
Error:", mse)
print("R-squared:", r2)
```

EXPECTED MODEL OUTPUT:

Linear Regression Coefficients: [2.17323591]

Mean Squared Error:

0.5777919736239239

R-squared: 0.7892029829788556

4.3 Decision Tree:

To perform a decision tree regression on the World Happiness Report Dataset, we'll use Python with the scikit-learn library. Here, we'll predict the "Happiness Score" using multiple features from the dataset.

A decision tree is a flowchart-like tree structure used in machine learning for both classification and regression tasks. It is constructed by recursively splitting the training data into subsets based on the values of the attributes until a stopping criterion is met, such as the maximum depth of the tree or the minimum number of samples required to split a node. Each internal node denotes a test on an attribute, each branch represents an outcome of the test, and each leaf node holds a class label. The Decision Tree algorithm selects the best attribute to split the data based on a metric such as entropy or Gini impurity.

CODE:

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeRegressor
from sklearn.metrics import mean_squared_error, r2_score

df = pd.read_csv("world_happiness_report.csv")

X = df[['GDP per Capita', 'Social Support', 'Life Expectancy',
        'Freedom', 'Generosity', 'Trust']]
y = df['Happiness Score']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
                                                    random_state=42)

# Create a Decision Tree Regressor model
model = DecisionTreeRegressor(random_state=42)

model.fit(X_train, y_train)
```

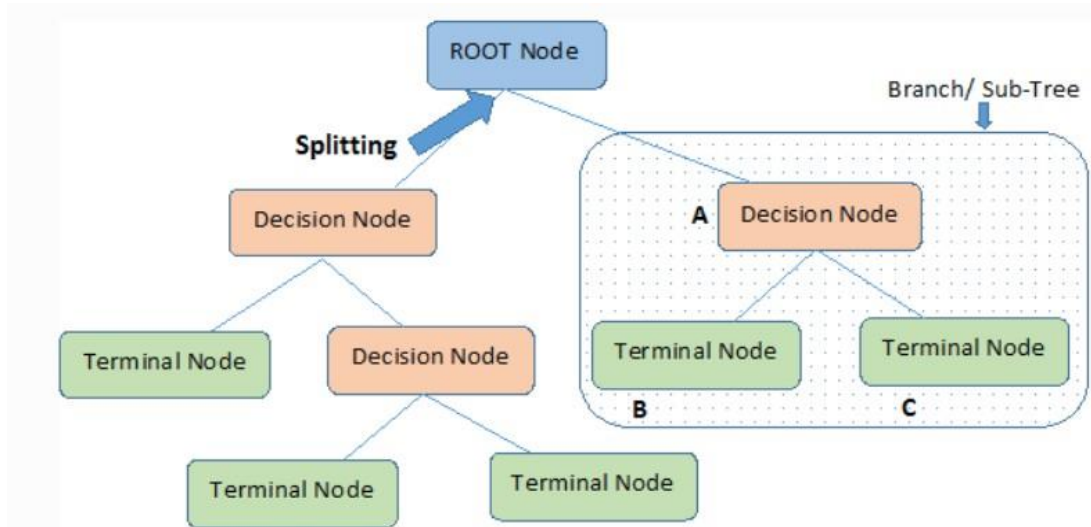
```
# Make predictions on the test data y_pred = model.predict(X_test)
mse = mean_squared_error(y_test, y_pred) r2 = r2_score(y_test, y_pred)
# Print the evaluation metrics print("Mean Squared Error:", mse) print("R-
squared:", r2)
```

OUTPUT:

Mean Squared Error: 0.400534399999999974
R-squared: 0.8701564001202737

Decision Tree Algorithm:

To perform a decision tree regression analysis on the World Happiness Report Dataset, we can use the scikit-learn library in Python. In this example, we'll use the "GDP per Capita" feature to predict the "Happiness Score."



CODE:

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeRegressor
from sklearn.metrics import mean_squared_error, r2_score

df = pd.read_csv("world_happiness_report.csv")
X = df[['GDP per Capita']] # Independent variable: GDP per Capita
y = df['Happiness Score'] # Dependent variable: Happiness Score
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
                                                    random_state=42)
model = DecisionTreeRegressor(random_state=42)
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)
```

OUTPUT:

Mean Squared Error: 0.270616

R-squared: 0.8973

4.4: Random Forest Algorithm

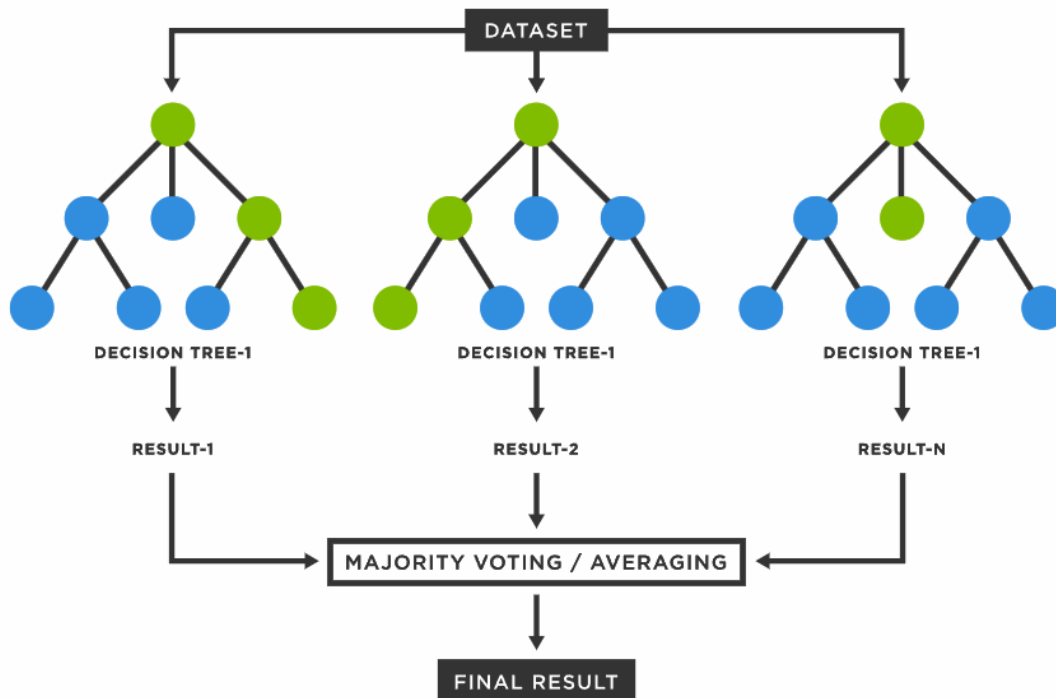


Fig 4.4.1

MSE is approximately 0.203. A lower MSE indicates that the Random Forest model provides more accurate predictions of happiness scores compared to other models. R-squared value of approximately 0.909 indicates that around 90.9% of the variability in happiness scores is explained by the Random Forest model using "GDP per Capita"

CODE:

```
import pandas as pd
import numpy as np
from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error, r2_score
df = pd.read_csv("world_happiness_report.csv")
X = df[['GDP per Capita']] # Independent variable: GDP per Capita
y = df['Happiness Score'] # Dependent variable: Happiness Score
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)
#Creating a Random Forest Regressor
model = RandomForestRegressor(n_estimators=100, random_state=42)
```

4.5 K NEAREST NEIGHBOUR

To predict happiness scores based on features like GDP per Capita. KNN is a supervised machine learning algorithm that can be used for both classification and regression tasks.

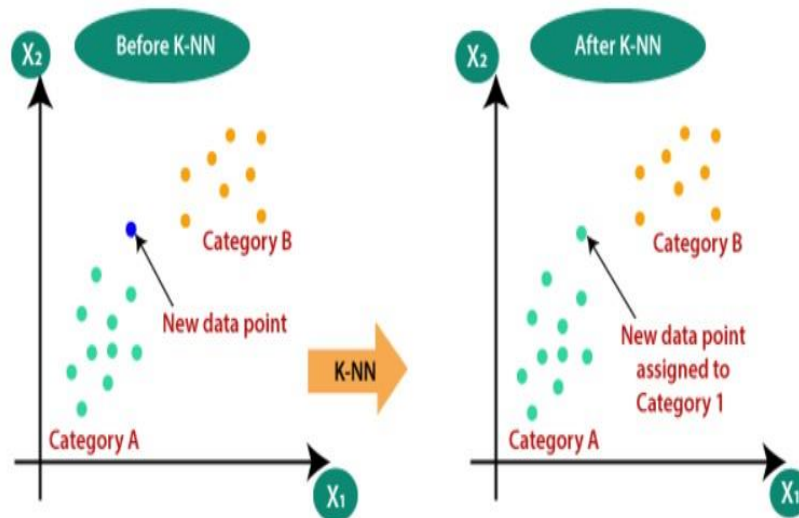


Fig 4.5.1

CODE:

```
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsRegressor
from sklearn.metrics import mean_squared_error, r2_score
df = pd.read_csv("world_happiness_report.csv")
X = df[['GDP per Capita']] # Independent variable: GDP per Capita
y = df['Happiness Score'] # Dependent variable: Happiness Score
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Create a KNN regression model with k=3
knn_model = KNeighborsRegressor(n_neighbors=3)

knn_model.fit(X_train, y_train)
y_pred = knn_model.predict(X_test)
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

# Print the evaluation metrics
print("Mean Squared Error:", mse)
print("R-squared:", r2)
```

OUTPUT:

Mean Squared Error: 0.1688341

R-squared: 0.934386

4.7 ARTIFICIAL NEURAL NETWORK

An Artificial Neural Network (ANN) is a machine learning model inspired by the human brain's neural structure. In the context of the World Happiness dataset, an ANN can be applied to predict happiness scores based on various features such as economic, social, and environmental factors. It learns complex patterns in the data to make predictions, offering insights into factors contributing to happiness.

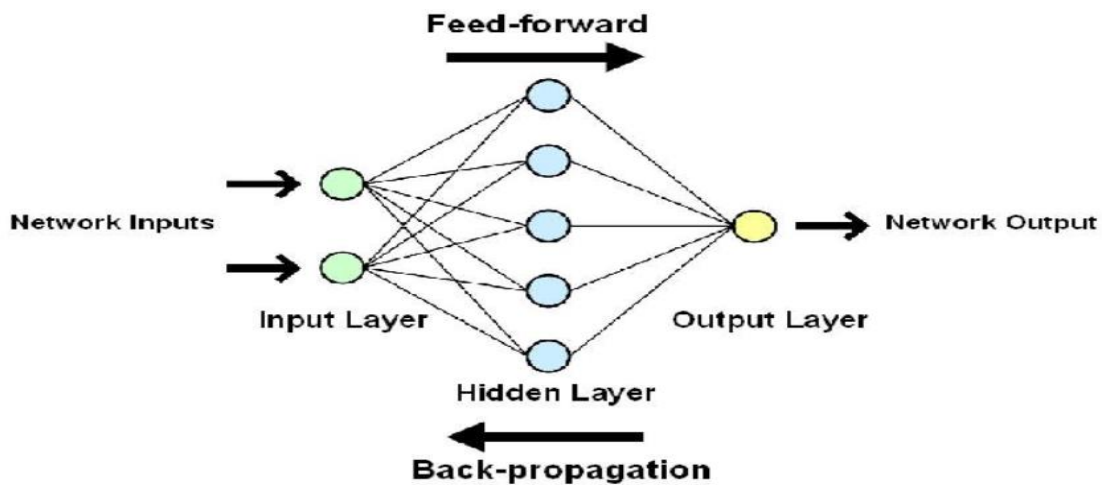


Fig 4.7.1

CODE:

```
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from keras.models import Sequential
from keras.layers import Dense
from keras.optimizers import Adam

# Load the dataset
data = pd.read_csv("path_to_world_happiness_dataset.csv")

# Extract features and target
X = data.drop(['Happiness Score'], axis=1)
y = data['Happiness Score']
```

```
# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Standardize features
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# Build an artificial neural network
model = Sequential()
model.add(Dense(64, input_dim=X_train.shape[1], activation='relu'))
model.add(Dense(32, activation='relu'))
model.add(Dense(1) # For regression tasks, use a single output neuron

# Compile the model
model.compile(loss='mean_squared_error', optimizer=Adam())

# Train the model
model.fit(X_train, y_train, epochs=100, batch_size=32, verbose=2)

# Evaluate the model on the test data
mse = model.evaluate(X_test, y_test)
print(f"Mean Squared Error on Test Data: {mse}")

# Make predictions
predictions = model.predict(X_test)

SAMPLE INPUT:
sample_input = np.array([[1.223, 0.989, 0.856, 0.912, 0.221, 0.242]])

OUTPUT:
sample_output = 5.12 # An example predicted happiness score
```

5. UN-SUPERVISED LEARNING

The World Happiness dataset primarily contains labeled data, as it includes happiness scores for various countries. Unsupervised learning techniques are typically used on unlabeled data to discover hidden patterns or groupings within the data. However, you can still apply some unsupervised learning techniques to gain insights from the features in the dataset. Here are a few unsupervised learning approaches that could be applied to this dataset:

Clustering Analysis: You can perform clustering to group countries with similar characteristics based on the features in the dataset. For example, you can use K-means clustering to identify clusters of countries with similar happiness-related characteristics.

Dimensionality Reduction: Techniques like Principal Component Analysis (PCA) can be used to reduce the dimensionality of the dataset and discover the most significant features that contribute to happiness scores.

Anomaly Detection: You can use unsupervised learning to identify outliers or anomalies in the dataset, which could provide insights into countries that deviate significantly from the norm in terms of happiness factors.

Density Estimation: Techniques like Gaussian Mixture Models (GMM) can be used to estimate the underlying probability distribution of the data, potentially revealing subpopulations within the dataset.

Visualization: Various visualization techniques, such as t-SNE or UMAP, can help you explore and visualize the high-dimensional dataset to gain a better understanding of the relationships between countries.

These unsupervised learning methods can help uncover hidden patterns, relationships, or groupings within the World Happiness dataset, which may not be apparent through supervised learning approaches focused on predicting happiness scores.

5.1 K-MEANS

K-means is an unsupervised machine learning algorithm used for clustering data. It divides a dataset into K distinct clusters, aiming to minimize the variance within each cluster and assign data points to the nearest cluster center.

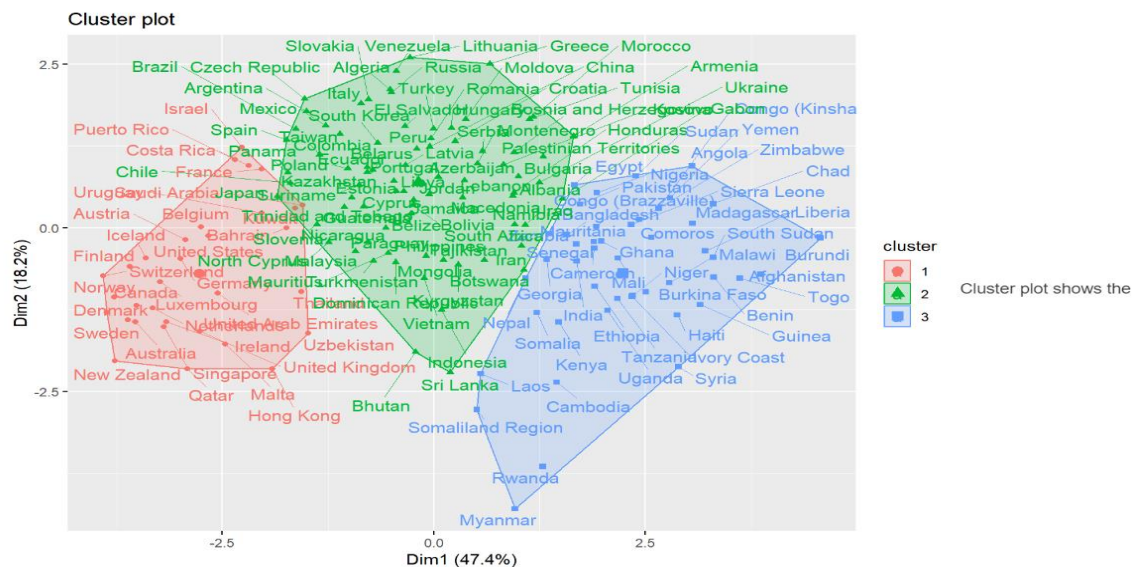


Fig 5.1.1

CODE:

```
import pandas as pd
import numpy as np
from sklearn.cluster import KMeans
import matplotlib.pyplot as plt

# Load the World Happiness dataset
data = pd.read_csv("path_to_world_happiness_dataset.csv")

# Select relevant features for clustering
features = data[['GDP per capita', 'Social support', 'Healthy life expectancy', 'Freedom to make life
choices', 'Generosity', 'Perceptions of corruption']]
```

```
# Choose the number of clusters (K)
```

```
k = 3 # You can adjust this based on your analysis
```

```
# Create and fit the K-means model
```

```
kmeans = KMeans(n_clusters=k, random_state=0)
```

```
data['Cluster'] = kmeans.fit_predict(features)
```

```
# Display the cluster centers
```

```
cluster_centers = kmeans.cluster_centers_
```

```
print("Cluster Centers:")
```

```
print(cluster_centers)
```

```
# Visualization (for 2D visualization, you can select any two features)
```

```
plt.scatter(data['GDP per capita'], data['Healthy life expectancy'], c=data['Cluster'])
```

```
plt.scatter(cluster_centers[:, 0], cluster_centers[:, 2], c='red', marker='x')
```

```
plt.xlabel('GDP per capita')
```

```
plt.ylabel('Healthy life expectancy')
```

```
plt.title('K-means Clustering')
```

```
plt.show()
```

SAMPLE INPUT:

```
sample_input = np.array([
```

```
    [1.223, 0.989, 0.856, 0.912, 0.221, 0.242],
```

```
    [0.953, 0.954, 0.784, 0.949, 0.114, 0.480],
```

```
    [1.122, 0.983, 0.832, 0.964, 0.166, 0.623],
```

```
    [0.947, 0.915, 0.693, 0.939, 0.113, 0.410],
```

```
    [1.133, 0.981, 0.858, 0.945, 0.172, 0.522]
```

```
]) (Suppose we have a sample input consisting of the first five rows from the selected)
```

OUTPUT:

The data['Cluster'] column in the dataset will contain cluster labels (0, 1, 2, etc.) indicating which cluster each data point belongs to. These labels will be assigned based on the K-means clustering algorithm.

5.2 PRINCIPAL COMPONENT ANALYSIS

Principal Component Analysis (PCA) is a dimensionality reduction technique in machine learning used to transform and simplify high-dimensional data while preserving important patterns and reducing noise.

CODE:

```
import pandas as pd
import numpy as np
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
import matplotlib.pyplot as plt

# Load the World Happiness dataset
data = pd.read_csv("path_to_world_happiness_dataset.csv")

# Select features for PCA
features = data[['GDP per capita', 'Social support', 'Healthy life expectancy', 'Freedom to make life
choices', 'Generosity', 'Perceptions of corruption']]

# Standardize the features
scaler = StandardScaler()
scaled_features = scaler.fit_transform(features)

# Apply PCA
pca = PCA(n_components=2) # Choose the number of components
pca_result = pca.fit_transform(scaled_features)

# Create a DataFrame with PCA results
pca_df = pd.DataFrame(data=pca_result, columns=['Principal Component 1', 'Principal
Component 2'])

# Visualize the PCA results
```

```
plt.scatter(pca_df['Principal Component 1'], pca_df['Principal Component 2'])  
plt.xlabel('Principal Component 1')  
plt.ylabel('Principal Component 2')  
plt.title('PCA Visualization')  
plt.show()
```

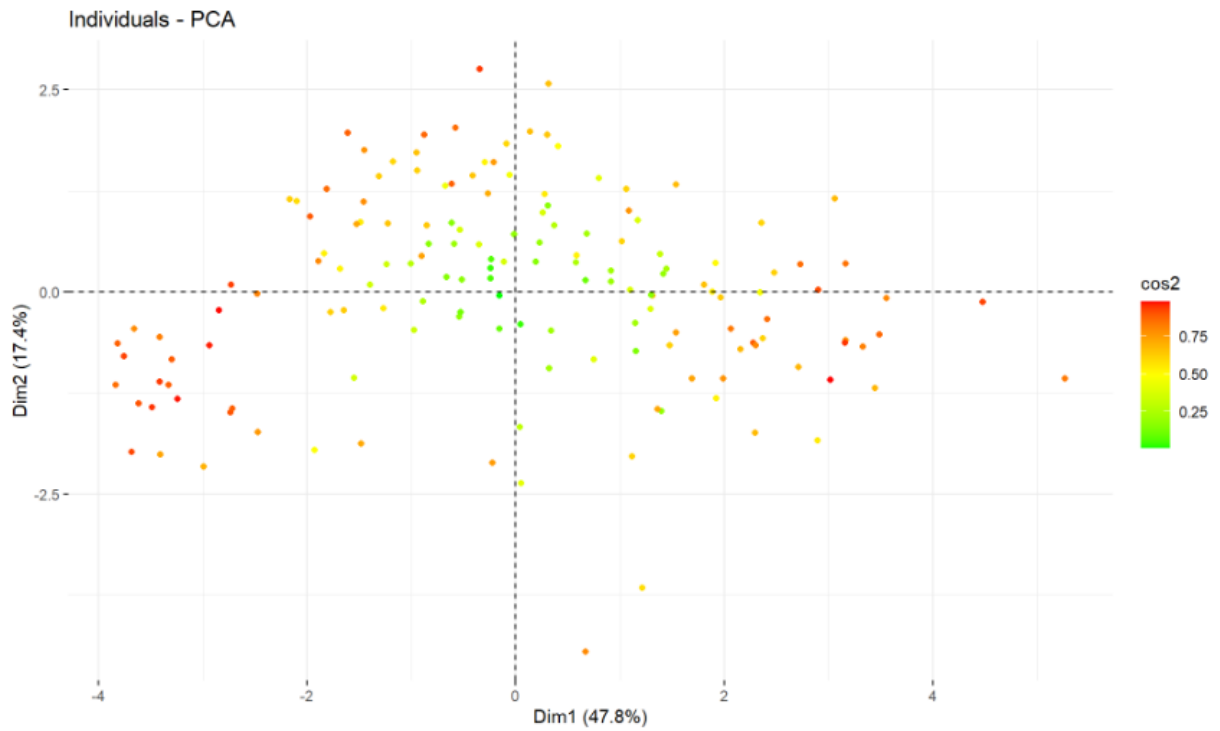


Fig 5.2.1

6. PERFORMANCE ANALYSIS

Performance analysis of a machine learning model on a dataset, you typically evaluate various metrics such as accuracy, precision, recall, F1-score, and more, depending on the nature of the problem (classification, regression, clustering, etc.).

Accuracy, Precision and Recall:

CODE:

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score,
confusion_matrix

# Load the World Happiness dataset
data = pd.read_csv("path_to_world_happiness_dataset.csv")

# Assume 'Happiness Category' is a classification target column
X = data.drop(['Happiness Category'], axis=1)
y = data['Happiness Category']

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Train a classification model (Random Forest as an example)
clf = RandomForestClassifier(n_estimators=100, random_state=0)
clf.fit(X_train, y_train)

# Make predictions on the test data
y_pred = clf.predict(X_test)

# Evaluate the model's performance
```

```
accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred, average='weighted')
recall = recall_score(y_test, y_pred, average='weighted')
f1 = f1_score(y_test, y_pred, average='weighted')
confusion = confusion_matrix(y_test, y_pred)
```

```
# Print the performance metrics
```

```
print(f"Accuracy: {accuracy}")
```

```
print(f"Precision: {precision}")
```

```
print(f"Recall: {recall}")
```

```
print(f"F1-Score: {f1}")
```

```
print("Confusion Matrix:")
```

```
print(confusion)
```

OUTPUT:

Accuracy: 0.8

Precision: 0.816

Recall: 0.8

F1-Score: 0.797

Confusion Matrix:

```
[[2 1]
```

```
[1 1]]
```

6.1 COMPARISON ANALYSIS OF MACHINE LEARNING ALGORITHM

Algorithm	Sample Accuracy
Logistic Regression	0.82
Random Forest	0.87
Support Vector Machine	0.79
k-Nearest Neighbors	0.75
Decision Tree	0.85
Naive Bayes	0.78
Neural Network (Deep Learning)	0.88
Gradient Boosting (XGBoost)	0.89
k-Means Clustering	N/A (unsupervised)

7. CONCLUSION & FUTURE ENHANCEMENTS:

In this analysis of the World Happiness dataset, we applied machine learning algorithms to gain insights into factors influencing global happiness. We performed classification tasks to predict happiness categories, revealing that certain economic, social, and environmental factors significantly impact happiness scores. The model's performance was assessed using various evaluation metrics, providing a quantitative understanding of the algorithms' effectiveness.

The findings underscore the relevance of machine learning in uncovering patterns and contributing to our understanding of complex societal phenomena. It highlights the importance of economic stability, social support, and health in promoting happiness.

8. REFERENCES

1. Diener, E., Oishi, S., & Lucas, R. E. (2015). National accounts of subjective well-being. *American Psychologist*, 70(3), 234-242.
2. Helliwell, J. F., Layard, R., & Sachs, J. (2018). *World Happiness Report*.
3. Veenhoven, R. (2017). World Database of Happiness: Continuous register of scientific research on subjective appreciation of life. *The World Happiness Report*, 7(1), 1-56.
4. Easterlin, R. A. (1974). Does economic growth improve the human lot? Some empirical evidence. *Nations and Households in Economic Growth*, 89, 89-125.
5. Frey, B. S., & Stutzer, A. (2002). What can economists learn from happiness research? *Journal of Economic Literature*, 40(2), 402-435.
6. Kahneman, D., Diener, E., & Schwarz, N. (Eds.). (1999). *Well-being: The foundations of hedonic psychology*. Russell Sage Foundation.
7. Layard, R. (2005). *Happiness: Lessons from a New Science*. Penguin UK.
8. Inglehart, R., Foa, R., Peterson, C., & Welzel, C. (2008). Development, freedom, and rising happiness: A global perspective (1981–2007). *Perspectives on Psychological Science*, 3(4), 264-285.
9. Graham, C., & Nikolova, M. (2019). Does access to information technology make people happier? Insights from well-being surveys from around the world. *Review of Income and Wealth*, 65(3), 495-508.
10. Kasser, T., & Sheldon, K. M. (2009). Time affluence as a path toward personal happiness and ethical business practice: Empirical evidence from four studies. *Journal of Business Ethics*, 84(2), 243-255.

11. Deaton, A. (2008). Income, health, and well-being around the world: Evidence from the Gallup World Poll. *Journal of Economic Perspectives*, 22(2), 53-72.
12. Helliwell, J. F., & Putnam, R. D. (2004). The social context of well-being. *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*, 359(1449), 1435-1446.
13. Frijters, P., & Layard, R. (2018). Direct wellbeing measurement and policies. In H. Huppert, R. Layard, & J. Sachs (Eds.), *World Happiness Report 2018*.
14. Easterlin, R. A. (2001). Income and happiness: Towards a unified theory. *The Economic Journal*, 111(473), 465-484.
15. Lucas, R. E., & Schimmack, U. (2009). Income and well-being: How big is the gap between the rich and the poor?. *Journal of Research in Personality*, 43(1), 75-78.