**PHASE-2:**

**Algorithm for Assessment of Marginal Workers in Tamil Nadu**

IMB_PROJECT

import pandas as pd import random

Define parameters for generating synthetic data

num_samples = 1000 # Number of data samples min_age = 18 max_age = 60 genders = ['Male', 'Female'] marital_statuses = ['Single', 'Married', 'Divorced', 'Widowed'] education_levels = ['No Education', 'Primary', 'Secondary', 'Higher Secondary', 'Graduate'] employment_statuses = ['Employed', 'Unemployed', 'Underemployed'] types_of_employment = ['Agriculture', 'Industry', 'Service', 'Others'] min_income = 1000 max_income = 50000 locations = ['Urban', 'Rural']

Create an empty DataFrame to store synthetic data

data = { 'Age': [], 'Gender': [], 'MaritalStatus': [], 'EducationLevel': [], 'EmploymentStatus': [], 'TypeOfEmployment': [], 'Income': [], 'Location': [], # Add more synthetic features as needed }

Generate synthetic data

for _ in range(num_samples): data['Age'].append(random.randint(min_age, max_age)) data['Gender'].append(random.choice(genders)) data['MaritalStatus'].append(random.choice(marital_statuses)) data['EducationLevel'].append(random.choice(education_levels)) data['EmploymentStatus'].append(random.choice(employment_statuses)) data['TypeOfEmployment'].append(random.choice(types_of_employment)) data['Income'].append(random.randint(min_income, max_income)) data['Location'].append(random.choice(locations)) # Add more synthetic features as needed
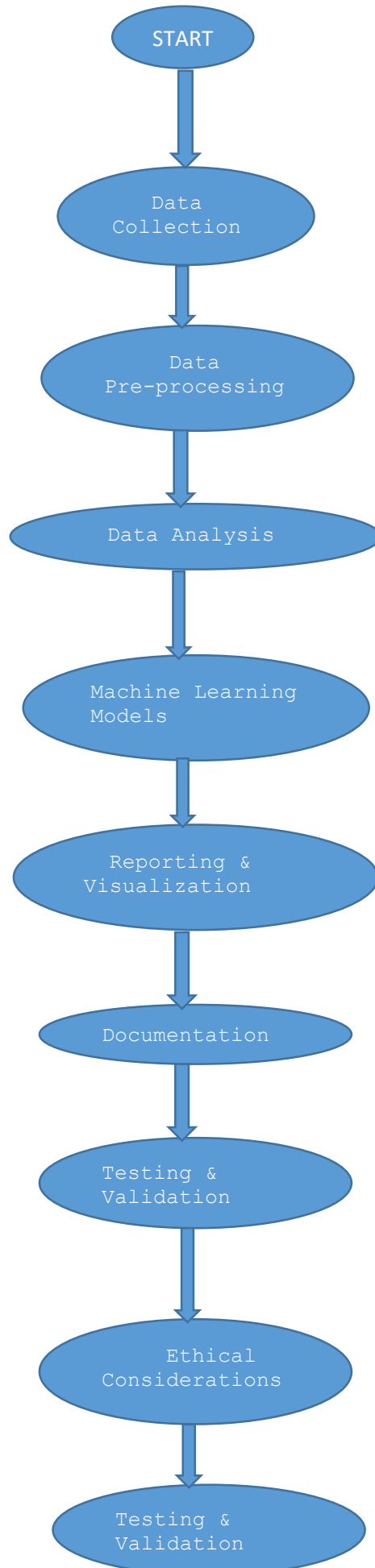
Create a Pandas DataFrame from the dictionary

df = pd.DataFrame(data)

Save the synthetic dataset to a CSV file

```
df.to_csv('synthetic_marginal_worker_data.csv', index=False)
```

**Flow Chart:**

START

↓

Data
Collection

↓

Data
Pre-processing

↓

Data Analysis

↓

Machine Learning
Models

↓

Reporting &
Visualization

↓

Documentation

↓

Testing &
Validation

↓

Ethical
Considerations

↓

Testing &
Validation

**Source Code:**

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

# Create synthetic data (you should replace this with real data)
data = {
    'Age': np.random.randint(18, 60, 1000),
    'EducationLevel': np.random.choice(['No Education', 'Primary',
'Secondary', 'Graduate'], 1000),
    'Income': np.random.randint(1000, 50000, 1000)
}

# Create a Pandas DataFrame
df = pd.DataFrame(data)

# Data preprocessing (feature encoding, handling missing values) would go
here if needed

# Split the data into features (X) and target (y)
X = df[['Age', 'EducationLevel']]
y = df['Income']

# Encode categorical variables if necessary

# Split 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)

# Create a Linear Regression model
model = LinearRegression()

# Train the model
model.fit(X_train, y_train)

# Make predictions
y_pred = model.predict(X_test)

# Evaluate the model
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

print(f"Mean Squared Error: {mse:.2f}")
print(f"R-squared (R2) Score: {r2:.2f}")
```

**Datasets of Assessment of Marginal Workers in Tamil Nadu**

1. **Demographic Information**:
    - Age
    - Gender
    - Marital status
    - Education level
    - Caste or community
2. **Employment Information**:
    - Employment status (e.g., employed, unemployed, underemployed)
    - Type of employment (e.g., agricultural, industrial, service)
    - Occupation
    - Sector of employment (e.g., informal, formal)
    - Income or wage levels
    - Duration of employment
3. **Geographic Information**:
    - Location (urban, rural)
    - District or region within Tamil Nadu
    - Accessibility to amenities (e.g., healthcare, education)
4. **Household Information**:
    - Household size
    - Household income
    - Dependency ratio (number of dependents per worker)
5. **Labor Force Participation**:
    - Participation in the labor force
    - Reasons for not participating (if not in the labor force)
    - Hours worked per day or week
6. **Migration Status**:
    - Internal or external migration (if applicable)
    - Reasons for migration
    - Duration of migration
7. **Social and Economic Indicators**:
    - Access to social services (e.g., healthcare, education)
    - Poverty status
    - Access to financial services
    - Access to government welfare programs
8. **Health and Well-being**:
    - Health status
    - Access to healthcare
    - Nutritional status
9. **Family and Household Dynamics**:
    - Family composition
    - Dependency on family members
    - Household assets
10. **Government Policies and Interventions**:
    - Government programs related to employment and social welfare
    - Participation in government schemes

11. **Socioeconomic Challenges and Barriers**:
   o Discrimination and social barriers faced by marginal workers
12. **Labor Market Trends**:
   o Employment trends in Tamil Nadu
   o Job opportunities and challenges
13. **Economic and Industrial Data**:
   o Economic indicators for Tamil Nadu
   o Data on key industries and sectors

Reference of the source link : https://tn.data.gov.in/catalog/marginal-workers-classified-age-industrial-category-and-sex-census-2011-india-and-states



## Linear Regression

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score
import matplotlib.pyplot as plt

# Load your dataset (replace 'your_dataset.csv' with the actual file path)
data = pd.read_csv(DDW_B06_3300_State_TAMIL_NADU-2011.csv')

# Define the features (independent variables) and the target variable
(Income)
X = data[['Feature1', 'Feature2', 'Feature3']]  # Add relevant features
y = data['Income']  # Replace 'Income' with your actual target variable
```

```python
# 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)

# Create a Linear Regression model
model = LinearRegression()

# Train the model on the training data
model.fit(X_train, y_train)

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

# Evaluate the model
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

print("Mean Squared Error:", mse)
print("R-squared (R2) Score:", r2)

# Visualize the model's predictions (optional)
plt.scatter(y_test, y_pred)
plt.xlabel("True Values")
plt.ylabel("Predictions")
plt.show()
```

1. **Load Data**: Load your dataset with features (independent variables) and the target variable (Income in this example).
2. **Define Features and Target**: Specify the features that you want to use for the linear regression model, as well as the target variable you want to predict.
3. **Split Data**: Split the data into training and testing sets to assess the model's performance.
4. **Create and Train the Model**: Create an instance of the Linear Regression model and train it on the training data.
5. **Make Predictions**: Use the trained model to make predictions on the test data.
6. **Evaluate the Model**: Calculate the Mean Squared Error (MSE) and R-squared (R2) score to assess the model's performance.
7. **Optional: Visualization**: You can visualize the model's predictions to see how well it fits the actual values.