the three stages of a reproducible workflow are:

Data cleaning and preprocessing: This stage involves preparing the data for analysis by checking for missing values, outliers, and inconsistencies. This stage aims to make the data ready for analysis by ensuring it is accurate, consistent, and complete. In this stage, you can do the following work:

Check for missing values and remove them or impute them.

Check for outliers and deal with them accordingly.

Check for inconsistencies in the data, such as spelling errors, and correct them.

Rename columns and select relevant features.

Convert data types if needed.

The folder structure for this stage can be organized as follows:

**Exploratory data analysis:** This stage involves exploring the data to gain insights and identify patterns. This stage aims to understand the data better and generate hypotheses for further analysis. In this stage, you can do the following work:

Compute summary statistics.

Create visualizations.

Identify patterns and correlations.

Check for assumptions.

## Generate hypotheses for further analysis.

The folder structure for this stage can be organized as follows:

analysis/

— exploratory/

| — data\_file\_1\_summary\_stats.csv

| — data\_file\_2\_summary\_stats.csv

| — data\_file\_1\_visualizations/

| — data\_file\_2\_visualizations/

— hypotheses/

— hypothesis 1.csv

**Modeling and inference:** This stage involves building models and testing hypotheses. This stage aims to draw conclusions from the data and make predictions. In this stage, you can do the following work:

Select appropriate models.

hypothesis\_2.csv

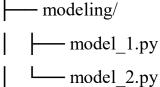
Test hypotheses.

Validate models.

Make predictions.

The folder structure for this stage can be organized as follows:

analysis/



```
inference/
results_1.csv
results 2.csv
```

## **Task 1: Reproducible workflow**

```
import pandas as pd
import os
# Set up folder structure
data_folder = 'data'
raw_data_folder = os.path.join(data_folder, 'raw_data')
cleaned data folder = os.path.join(data folder, 'cleaned data')
analysis folder = 'analysis'
exploratory_folder = os.path.join(analysis_folder, 'exploratory')
hypotheses folder = os.path.join(analysis folder, 'hypotheses')
if not os.path.exists(raw data folder):
  os.makedirs(raw data folder)
if not os.path.exists(cleaned_data_folder):
  os.makedirs(cleaned_data_folder)
if not os.path.exists(exploratory folder):
  os.makedirs(exploratory_folder)
if not os.path.exists(hypotheses folder):
  os.makedirs(hypotheses_folder)
```

```
# Load raw data
raw data = pd.read csv(os.path.join(raw data folder, 'my data.csv'))
# Data cleaning and preprocessing
cleaned_data = raw_data.dropna() # remove missing values
# Save cleaned data
cleaned data.to csv(os.path.join(cleaned data folder, 'my data cleaned.csv'))
# Exploratory data analysis
summary_stats = cleaned_data.describe() # compute summary statistics
summary stats.to csv(os.path.join(exploratory folder,
'my_data_summary_stats.csv'))
import matplotlib.pyplot as plt
import seaborn as sns
# Create visualizations
sns.scatterplot(data=cleaned data, x='study time', y='final grade')
plt.savefig(os.path.join(exploratory folder, 'my data visualizations',
'scatterplot.png'))
sns.histplot(data=cleaned_data, x='final_grade')
```

```
plt.savefig(os.path.join(exploratory_folder, 'my_data_visualizations',
'histogram.png'))
sns.barplot(data=cleaned_data, x='gender', y='final_grade')
plt.savefig(os.path.join(exploratory folder, 'my data visualizations',
'barplot.png'))
sns.boxplot(data=cleaned_data, y='final_grade')
plt.savefig(os.path.join(exploratory_folder, 'my_data_visualizations',
'boxplot.png'))
# Modeling and inference
from sklearn.linear model import LinearRegression
# Fit linear regression model
model = LinearRegression()
model.fit(cleaned data[['study time']], cleaned data['final grade'])
# Make predictions
predictions = model.predict(cleaned_data[['study_time']])
results = pd.DataFrame({'actual': cleaned data['final grade'], 'predicted':
predictions})
results.to csv(os.path.join(hypotheses folder, 'my hypothesis results.csv'))
```