**app.py**

```python
1   import streamlit as st
2   import pandas as pd
3   import os
4   import sys
5
6   # Add the current directory to path to allow imports from src
7   sys.path.append(os.getcwd())
8
9   from src.eda import load_data, show_stats, plot_correlation, plot_distribution,
    show_missing_values, plot_pairplot, impute_missing_values, impute_all_missing_values,
    convert_to_numeric
10  from src.model import train_linear_regression, train_polynomial_regression, evaluate_model,
    plot_regression_results, train_knn_regression, train_random_forest_regression,
    generate_model_explanation, plot_actual_vs_predicted
11
12  # Set page config
13  st.set_page_config(page_title="EDA & Regression Analysis", layout="wide")
14
15  st.title("Nuclear Energy Insights Dashboard & EDA")
16
17  # --- Sidebar: Data Loading ---
18  st.sidebar.header("1. Upload Data")
19  uploaded_file = st.sidebar.file_uploader("Upload your CSV file", type=["csv"])
20
21  # Default file path
22  default_file_path = "archive/us_nuclear_generating_statistics_1971_2021.csv"
23
24  # Initialize session state for dataframe if not exists
25  if 'df' not in st.session_state:
26      st.session_state.df = None
27
28  # Load data logic
29  if uploaded_file is not None:
30      # Check if we need to reload (e.g. new file uploaded)
31      # Simple check: just reload. For optimization, could check file name.
32      # For now, if uploaded_file changes, Streamlit re-runs script, so we reload.
33      st.session_state.df = load_data(uploaded_file)
34      st.sidebar.success("File uploaded successfully!")
35  elif st.session_state.df is None and os.path.exists(default_file_path):
36      st.sidebar.info(f"Using default dataset: {os.path.basename(default_file_path)}")
37      st.session_state.df = load_data(default_file_path)
38  elif st.session_state.df is None:
39      st.sidebar.warning("Please upload a CSV file to proceed.")
40
41  # --- Main App Logic ---
42  if st.session_state.df is not None:
43      df = st.session_state.df # Local alias for convenience
44      # Sidebar Navigation
```

```python
      page = st.sidebar.radio("Navigate", ["Exploratory Data Analysis (EDA)", "Regression
   Modeling"])

      if page == "Exploratory Data Analysis (EDA)":
          st.header("🔍 Exploratory Data Analysis")

          # Data Overview
          st.subheader("Dataset Overview")
          st.write(f"Shape: {df.shape[0]} rows, {df.shape[1]} columns")
          st.dataframe(df.head())

          # Stats
          st.subheader("Descriptive Statistics")
          st.write(df.describe())

          # Missing Values
          st.subheader("Missing Values")
          missing_vals = show_missing_values(df)
          st.write(missing_vals)

          # Imputation
          if missing_vals.sum() > 0:
              st.markdown("### Impute Missing Values")
              cols_with_missing = missing_vals[missing_vals > 0].index.tolist()

              if cols_with_missing:
                  c1, c2, c3 = st.columns([2, 1, 1])
                  with c1:
                      col_to_impute = st.selectbox("Select Column to Impute",
   cols_with_missing)
                  with c2:
                      imp_strategy = st.selectbox("Strategy", ["Mean", "Median", "Mode"])
                  with c3:
                      if st.button("Apply Imputation"):
                          st.session_state.df = impute_missing_values(st.session_state.df,
   col_to_impute, imp_strategy)
                          st.success(f"Imputed {col_to_impute} with {imp_strategy}")
                          st.rerun()

              st.markdown("#### Bulk Imputation")
              c_bulk1, c_bulk2 = st.columns([2, 1])
              with c_bulk1:
                  bulk_strategy = st.selectbox("Bulk Strategy (All Columns)", ["Mean",
   "Median", "Mode"])
              with c_bulk2:
                  if st.button("Impute All"):
                      st.session_state.df = impute_all_missing_values(st.session_state.df,
   bulk_strategy)
                      st.success(f"Imputed all valid columns with {bulk_strategy}")
                      st.rerun()
```

```
 91            # Visualizations
 92            st.subheader("Visualizations")
 93
 94            col1, col2 = st.columns(2)
 95
 96            with col1:
 97                st.markdown("### Correlation Heatmap")
 98                fig_corr = plot_correlation(df)
 99                if fig_corr:
100                    st.pyplot(fig_corr)
101
102            with col2:
103                st.markdown("### Distribution Plot")
104                numeric_cols = df.select_dtypes(include=['float64', 'int64']).columns.tolist()
105                if numeric_cols:
106                    selected_col = st.selectbox("Select column for distribution", numeric_cols)
107                    fig_dist = plot_distribution(df, selected_col)
108                    st.pyplot(fig_dist)
109                else:
110                    st.write("No numeric colums for distribution plot.")
111
112      elif page == "Regression Modeling":
113            st.header("📈 Regression Modeling")
114
115            # Column Selection
116            numeric_cols = df.select_dtypes(include=['float64', 'int64']).columns.tolist()
117
118            # Check if we have enough numeric columns; if not, try to convert
119            if len(numeric_cols) < 2:
120                with st.spinner("Attempting to convert text columns to numbers..."):
121                    st.session_state.df = convert_to_numeric(st.session_state.df)
122                    df = st.session_state.df # Refresh local alias
123                    numeric_cols = df.select_dtypes(include=['float64',
     'int64']).columns.tolist()
124
125                    if len(numeric_cols) >= 2:
126                        st.success(f"Successfully converted data! Found {len(numeric_cols)} numeric
     columns.")
127
128            if len(numeric_cols) < 2:
129                st.error("Dataset needs at least 2 numeric columns for regression.")
130                st.write("Current Numeric Columns:", numeric_cols)
131                st.write("All Columns & Types:", df.dtypes)
132            else:
133                col1, col2 = st.columns(2)
134                with col1:
135                    target_col = st.selectbox("Select Target Variable (Y)", numeric_cols,
     index=len(numeric_cols)-1)
136                with col2:
137                    feature_options = [c for c in numeric_cols if c != target_col]
138                    # Auto-select the first feature by default to avoid empty state error
```

```python
                    default_feat = [feature_options[0]] if feature_options else None
                    feature_col = st.multiselect("Select Feature Variable(s) (X)",
        feature_options, default=default_feat)

                if not feature_col:
                    st.warning("Please select at least one feature variable.")
                else:
                    model_type = st.radio("Select Model Type", ["Linear Regression", "Polynomial
        Regression", "KNN Regression", "Random Forest Regression"])

                    degree = 2
                    k_neighbors = 5
                    n_estimators = 100

                    if model_type == "Polynomial Regression":
                        degree = st.slider("Select Polynomial Degree", 2, 5, 2)
                    elif model_type == "KNN Regression":
                        k_neighbors = st.slider("Select K Neighbors", 1, 20, 5)
                    elif model_type == "Random Forest Regression":
                        n_estimators = st.slider("Select Number of Trees (Estimators)", 10, 500,
        100, step=10)

                    if st.button("Train Model"):
                        # Create a subset for training
                        train_df = df.dropna(subset=feature_col + [target_col])

                        if len(train_df) == 0:
                            st.error("No data left after removing missing values. Please check
        your data.")
                        else:
                            if len(df) != len(train_df):
                                st.warning(f"Dropped {len(df) - len(train_df)} rows containing
        missing values.")

                            X = train_df[feature_col].values # Multiselect returns list, so this
        works for both single and multi
                            y = train_df[target_col].values

                            # Ensure X is 2D
                            if len(X.shape) == 1:
                                X = X.reshape(-1, 1)

                            poly_features = None # Default

                            if model_type == "Linear Regression":
                                model = train_linear_regression(X, y)
                            elif model_type == "Polynomial Regression":
                                model, poly_features = train_polynomial_regression(X, y, degree)
                            elif model_type == "KNN Regression":
                                model = train_knn_regression(X, y, k_neighbors)
                            elif model_type == "Random Forest Regression":
```

```
184                    model = train_random_forest_regression(X, y, n_estimators)
185
186                    # Evaluate happens here for all because logic is shared except for poly
       transform
187                    mse, r2, y_pred = evaluate_model(model, X, y, poly_features) #
       evaluate_model generates predictions too
188
189                    # metrics
190                    st.success("Model Trained!")
191                    m_col1, m_col2 = st.columns(2)
192                    m_col1.metric("R2 Score", f"{r2:.4f}")
193                    m_col2.metric("MSE", f"{mse:.4f}")
194
195                    with st.expander("ℹ️ How to interpret these results?"):
196                        st.write("""
197                        **1. R2 Score (0 to 1):**
198                        - Represents accuracy. **1.0 (100%)** is perfect.
199                        - **< 0.3**: Weak prediction.
200                        - **0.3 - 0.7**: Moderate.
201                        - **> 0.7**: Strong.
202
203                        **2. Mean Squared Error (MSE):**
204                        - The average squared difference between actual and predicted values.
205                        - **Lower is better**. 0 means no error.
206
207                        **3. Regression Plot:**
208                        - **Blue Dots**: The model's predictions.
209                        - **Red Line**: Perfect prediction (Actual = Predicted).
210                        - **Goal**: Points should be as close to the red line as possible.
211                        """)
212
213                    # Plot
214                    st.subheader("Regression Plot")
215                    if len(feature_col) > 1:
216                        # Multi-feature: Plot Actual vs Predicted
217                        fig_reg = plot_actual_vs_predicted(y, y_pred, title=f"{model_type}
       (Actual vs Predicted)")
218                        st.pyplot(fig_reg)
219                        st.info("Note: When using multiple features, we plot 'Actual vs
       Predicted' because we cannot easily visualize >3 dimensions.")
220                    else:
221                        # Single feature: Standard regression plot
222                        title = f"{model_type}"
223                        if model_type == "Polynomial Regression":
224                            title += f" (Degree: {degree})"
225                        elif model_type == "KNN Regression":
226                            title += f" (K: {k_neighbors})"
227                        elif model_type == "Random Forest Regression":
228                            title += f" (Trees: {n_estimators})"
229
230                        fig_reg = plot_regression_results(X, y, y_pred, title=title)
```

```
231                        st.pyplot(fig_reg)

232

233                    # Explanation
234                    st.subheader("Model Insights")
235                    explanation = generate_model_explanation(model, model_type, feature_col,
       target_col)
236                    st.markdown(explanation)

237

238            # Suggestion for non-linear check
239            st.info("Tip: If R2 is low for Linear Regression, try Polynomial Regression to
       capture non-linear relationships.")

240
```