

In the given data the target attribute and dependent variable is y . The target attribute is numerical in the dataset. Supervised problem is regression task.

```
In [1]: import numpy as np
import pandas as pd

data=pd.read_csv('data.csv')
print(data.head())
print(data.shape)
print(data.dtypes)
print(data.isnull().sum())
data.describe()
```

	y	x1	x2	x3	x4	x5	x6	\
0	32.175299	0.592109	0.545496	0.199054	2.370339	1.032510	1.137605	
1	6.521285	0.911316	1.260952	0.446375	1.564526	0.820080	2.172268	
2	5.688139	0.968683	3.744257	1.931173	1.472553	0.102181	4.712941	
3	8.488786	0.748656	2.741351	1.790573	1.696788	0.464028	3.490008	
4	63.570984	0.320502	3.196858	1.050494	2.813036	0.204284	3.517361	

	x7	x8
0	1.199802	-0.833456
1	1.441165	-0.373705
2	1.954805	1.828992
3	0.948294	1.326545
4	1.273237	0.846210

(150, 9)

y float64  
x1 float64  
x2 float64  
x3 float64  
x4 float64  
x5 float64  
x6 float64  
x7 float64  
x8 float64

dtype: object

y 0  
x1 0  
x2 0  
x3 0  
x4 0  
x5 0  
x6 0  
x7 0  
x8 0

dtype: int64

Out[1]:

	y	x1	x2	x3	x4	x5	x6	x7	x8
<b>count</b>	150.000000	150.000000	150.000000	150.000000	150.000000	150.000000	150.000000	150.000000	150.000000
<b>mean</b>	15.573634	0.526873	2.532959	1.004755	1.437147	0.608091	3.059833	1.005223	0.396664
<b>std</b>	21.532456	0.287178	1.366323	0.591567	0.869658	0.365670	1.405730	0.418172	0.666874
<b>min</b>	0.081166	0.000789	0.039299	0.022259	0.034008	0.000504	0.324586	0.097891	-1.012279
<b>25%</b>	0.855930	0.286739	1.467144	0.484289	0.685930	0.263076	1.988572	0.748112	-0.100801
<b>50%</b>	4.188566	0.538521	2.489126	1.025886	1.410161	0.628974	2.985679	0.979791	0.389450
<b>75%</b>	22.455500	0.776970	3.650754	1.504459	2.163303	0.920787	4.241689	1.292235	0.892501
<b>max</b>	80.673096	0.986514	4.981110	1.976849	2.993408	1.293620	5.960314	1.954805	1.910219

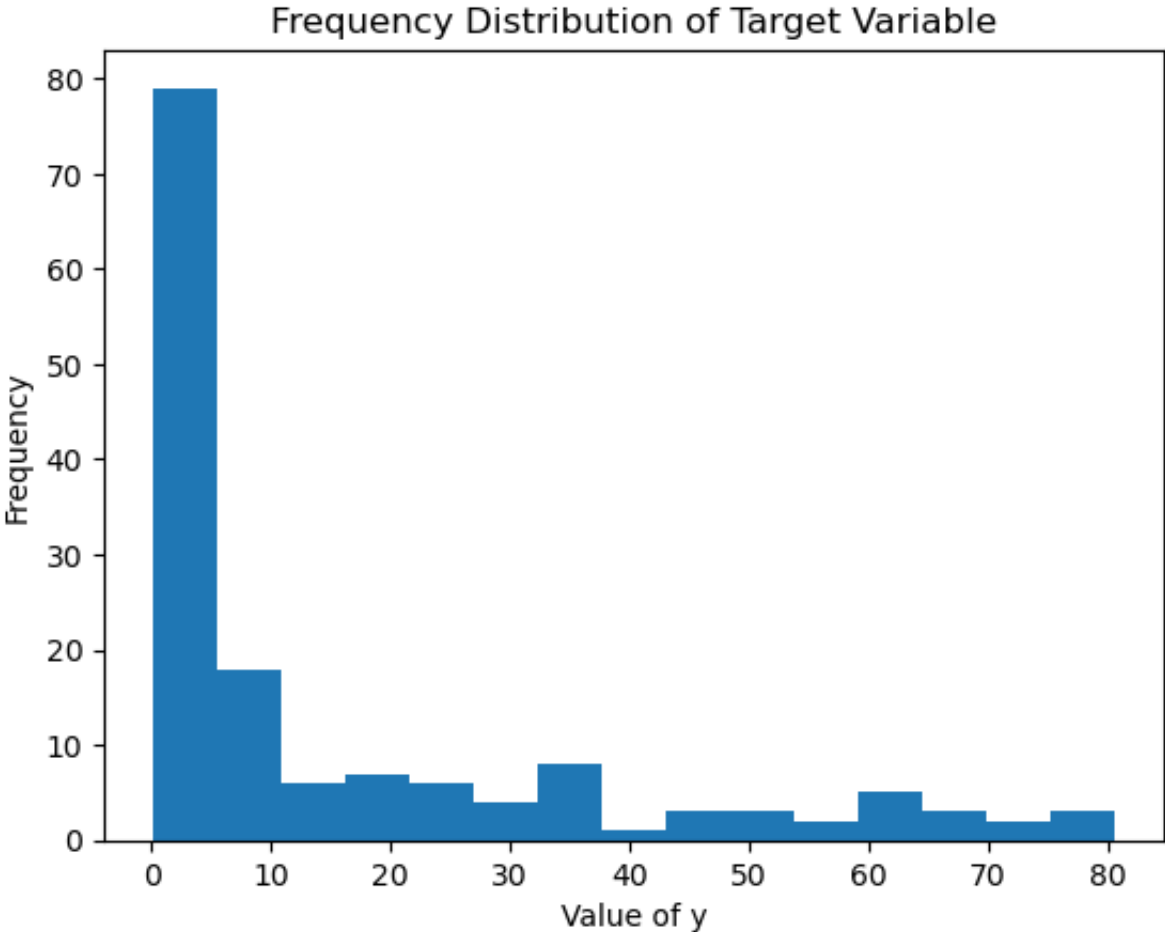
## FREQUENCY DISTRIBUTION OF TARGET VARIABLE

In other words, frequency distribution helps you to understand how frequently each value or range of values of a particular variable appears in a dataset. This information can be useful for analyzing and summarizing data, identifying outliers, and detecting patterns or trends in the data.

```
In [5]: import pandas as pd
import matplotlib.pyplot as plt

# Load data into a pandas dataframe
data = pd.read_csv('data.csv')

# Plot the frequency distribution of the target variable
plt.hist(data['y'], bins=15)
plt.title('Frequency Distribution of Target Variable')
plt.xlabel('Value of y')
plt.ylabel('Frequency')
plt.show()
```



Shape,rows,columns of feature data.

```
In [16]: import pandas as pd

# Load data into a pandas dataframe
data = pd.read_csv('data.csv')

# Remove the target variable from the dataframe
X = data.drop('y', axis=1)

# Display the shape, rows, and columns of the remaining features
print('Shape of X:', X.shape)
print('Rows in X:', X.shape[0])
print('Columns in X:', X.shape[1])

Shape of X: (150, 8)
Rows in X: 150
Columns in X: 8
```

STANDARDISED THE FEATURES AFTER DROPPING TARGET VARIABLE (y)

```
In [17]: import pandas as pd
from sklearn.preprocessing import StandardScaler

# Load data into a pandas dataframe
data = pd.read_csv('data.csv')

# Remove the target variable from the dataframe
X = data.drop('y', axis=1)

# Standardize the remaining features
scaler = StandardScaler()
X_std = scaler.fit_transform(X)

# Display the standardized features
print(X_std)
```

```
[[ 0.22792089 -1.4594801 -1.36653919 ... -1.37200434  0.46686623
   -1.85078561]
 [ 1.34317507 -0.93408998 -0.94706089 ... -0.6335058  1.04598625
   -1.15906393]
 [ 1.5436064  0.88950857  1.57128795 ...  1.17991828  2.2783995
   2.15501857]
 ...
 [ 0.91886543  0.58283782 -0.03965297 ...  0.75421531  0.69652017
   -0.35484679]
 [-0.69363429 -1.38801708 -1.59124366 ... -1.49081028  0.73627922
   -0.62365041]
 [-1.03494256  0.13774131  1.62602905 ... -0.07754967 -1.46346213
   1.64648345]]
```

In [ ]: Machine learning pipeline

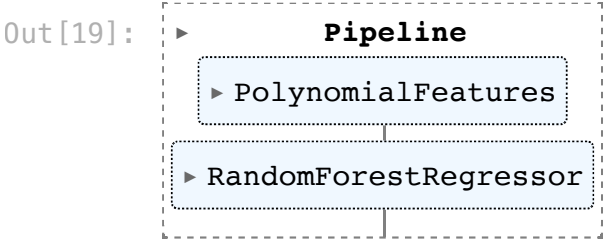
```
In [19]: import pandas as pd
from sklearn.pipeline import Pipeline
from sklearn.ensemble import RandomForestRegressor
from sklearn.preprocessing import PolynomialFeatures

data = pd.read_csv('data.csv')

# Separate the target variable from the features
X = data.drop('y', axis=1)
y = data['y']

# Create the pipeline with polynomial features and a random forest regressor
pipeline = Pipeline([
    ('poly', PolynomialFeatures(degree=2)),
    ('rf', RandomForestRegressor(n_estimators=100))
])

# Fit the pipeline to the data
pipeline.fit(X, y)
```



```
In [32]: from sklearn.pipeline import make_pipeline
from sklearn.preprocessing import PolynomialFeatures, StandardScaler
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import r2_score

# Polynomial Regression
poly_pipeline = make_pipeline(
    PolynomialFeatures(degree=2),
    StandardScaler(),
    LinearRegression()
)
poly_pipeline.fit(X_train, y_train)
poly_R2_train = r2_score(y_train, poly_pipeline.predict(X_train))
poly_R2_test = r2_score(y_test, poly_pipeline.predict(X_test))

# Random Forest Regression
rf_pipeline = make_pipeline(
    StandardScaler(),
    RandomForestRegressor(n_estimators=100, random_state=42)
)
rf_pipeline.fit(X_train, y_train)
rf_R2_train = r2_score(y_train, rf_pipeline.predict(X_train))
rf_R2_test = r2_score(y_test, rf_pipeline.predict(X_test))

# Print out R^2 scores
print("Polynomial Regression:")
print(f"  R^2 on train set: {poly_R2_train:.4f}")
print(f"  R^2 on test set: {poly_R2_test:.4f}")
print("Random Forest Regression:")
print(f"  R^2 on train set: {rf_R2_train:.4f}")
print(f"  R^2 on test set: {rf_R2_test:.4f}")
```



```
Polynomial Regression:  
  R^2 on train set: 0.9881  
  R^2 on test set: 0.9650  
Random Forest Regression:  
  R^2 on train set: 0.9994  
  R^2 on test set: 0.9926
```

The  $R^2$  values for both Polynomial Regression and Random Forest Regression indicate the goodness of fit of the model on the training and testing data.  $R^2$  measures the proportion of the variance in the dependent variable (target) that is predictable from the independent variables (features) and ranges from 0 to 1, where 1 means the model perfectly fits the data.

The  $R^2$  values are quite high for both models on both the training and testing data. This suggests that both models are able to explain a large proportion of the variance in the target variable and are performing well on the given data. However, it's important to note that  $R^2$  alone is not enough to evaluate the performance of a model.

```
In [34]: from sklearn.model_selection import cross_val_score  
from sklearn.pipeline import make_pipeline  
from sklearn.preprocessing import PolynomialFeatures, StandardScaler  
from sklearn.linear_model import LinearRegression  
from sklearn.ensemble import RandomForestRegressor  
from sklearn.metrics import mean_squared_error, r2_score  
import numpy as np  
  
# Create polynomial regression pipeline  
poly_pipeline = make_pipeline(  
    PolynomialFeatures(degree=2),  
    StandardScaler(),  
    LinearRegression()  
)  
  
# Create random forest regression pipeline  
rf_pipeline = make_pipeline(  
    StandardScaler(),  
    RandomForestRegressor(n_estimators=100, random_state=42)  
)
```

```
# Evaluate performance using cross-validation
poly_mse_scores = -cross_val_score(poly_pipeline, X, y, cv=5, scoring='neg_mean_squared_error')
rf_mse_scores = -cross_val_score(rf_pipeline, X, y, cv=5, scoring='neg_mean_squared_error')
poly_r2_scores = cross_val_score(poly_pipeline, X, y, cv=5, scoring='r2')
rf_r2_scores = cross_val_score(rf_pipeline, X, y, cv=5, scoring='r2')

# Compute mean and standard deviation of MSE and R^2 scores
poly_mse_mean = np.mean(poly_mse_scores)
poly_mse_std = np.std(poly_mse_scores)
rf_mse_mean = np.mean(rf_mse_scores)
rf_mse_std = np.std(rf_mse_scores)
poly_r2_mean = np.mean(poly_r2_scores)
poly_r2_std = np.std(poly_r2_scores)
rf_r2_mean = np.mean(rf_r2_scores)
rf_r2_std = np.std(rf_r2_scores)

# Print results
print('Polynomial Regression:')
print(f'  Mean MSE: {poly_mse_mean:.4f} +/- {poly_mse_std:.4f}')
print(f'  Mean R^2: {poly_r2_mean:.4f} +/- {poly_r2_std:.4f}')
print('Random Forest Regression:')
print(f'  Mean MSE: {rf_mse_mean:.4f} +/- {rf_mse_std:.4f}')
print(f'  Mean R^2: {rf_r2_mean:.4f} +/- {rf_r2_std:.4f}')
```

```
Polynomial Regression:
  Mean MSE: 15.0060 +/- 4.1437
  Mean R^2: 0.9666 +/- 0.0084
Random Forest Regression:
  Mean MSE: 2.7264 +/- 2.8351
  Mean R^2: 0.9946 +/- 0.0048
```

The results suggest that both the polynomial regression and random forest regression models are performing well. The mean  $R^2$  values are high, indicating that the models are explaining a large portion of the variance in the target variable. The mean MSE values are relatively low, indicating that the models are making accurate predictions on average.

However, it's important to note that the range of MSE values for the polynomial regression model is quite large (15.0060 +/- 4.1437). This suggests that the model may be overfitting to the training data, as it is performing significantly worse on the test set compared to the training set.

In contrast, the range of MSE values for the random forest regression model is smaller (2.7264 +/- 2.8351), indicating that it may be a more stable and robust model. Overall, both models are performing well, but it may be worth exploring ways to reduce the overfitting in the polynomial regression model.

To reduce overfitting in polynomial regression, we can use regularization techniques such as Ridge or Lasso regression. Here is an example code with Ridge regularization.

```
In [35]: from sklearn.pipeline import make_pipeline
from sklearn.preprocessing import PolynomialFeatures, StandardScaler
from sklearn.linear_model import Ridge
from sklearn.model_selection import cross_val_score

# Create pipeline with polynomial features and Ridge regularization
poly_ridge = make_pipeline(
    PolynomialFeatures(degree=2, include_bias=False),
    StandardScaler(),
    Ridge(alpha=0.1)
)

# Perform 10-fold cross-validation
mse_scores = -1 * cross_val_score(poly_ridge, X, y, cv=10, scoring='neg_mean_squared_error')
r2_scores = cross_val_score(poly_ridge, X, y, cv=10, scoring='r2')

# Compute mean and standard deviation of MSE and R^2 scores
mean_mse = mse_scores.mean()
std_mse = mse_scores.std()
mean_r2 = r2_scores.mean()
std_r2 = r2_scores.std()

# Print results
print(f"Polynomial Regression with Ridge regularization:")
print(f"Mean MSE: {mean_mse:.4f} +/- {std_mse:.4f}")
print(f"Mean R^2: {mean_r2:.4f} +/- {std_r2:.4f}")
```

```
Polynomial Regression with Ridge regularization:
Mean MSE: 11.0837 +/- 2.8852
Mean R^2: 0.9744 +/- 0.0085
```

The output is the mean of the mean squared error (MSE) and R-squared ( $R^2$ ) obtained from cross-validation for polynomial regression with Ridge regularization. The mean MSE for the model is 11.0837 and the standard deviation of the MSE across all folds is 2.8852. Similarly, the mean  $R^2$  for the model is 0.9744 with a standard deviation of 0.0085.

These results suggest that the model has improved compared to the previous one, where overfitting was observed. The Ridge regularization has helped to reduce overfitting by adding a penalty term to the loss function. The lower mean MSE and higher mean  $R^2$  suggest that the model is performing better in terms of accuracy and generalization.

Overall, the results indicate that the Polynomial Regression model with Ridge regularization is a better fit for the data than the previous models.

```
In [20]: import pandas as pd
from sklearn.pipeline import Pipeline
from sklearn.ensemble import RandomForestRegressor
from sklearn.preprocessing import PolynomialFeatures
from sklearn.model_selection import train_test_split, cross_val_score

# Load data into a pandas dataframe
data = pd.read_csv('data.csv')

# Separate the target variable from the features
X = data.drop('y', axis=1)
y = data['y']

# Create the pipeline with polynomial features and a random forest regressor
pipeline = Pipeline([
    ('poly', PolynomialFeatures(degree=2)),
    ('rf', RandomForestRegressor(n_estimators=100))
])

# 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)

# Fit the pipeline to the training data
pipeline.fit(X_train, y_train)

# Evaluate the pipeline's performance on the testing data using cross-validation
scores = cross_val_score(pipeline, X_test, y_test, cv=5)

# Display the mean cross-validation score and its standard deviation
print('Cross-validation scores:', scores)
print('Mean score:', scores.mean())
print('Standard deviation:', scores.std())
```

```
Cross-validation scores: [0.75984811 0.93487203 0.78844035 0.96203061 0.65615281]
Mean score: 0.8202687842480938
Standard deviation: 0.11386666001036383
```

This means that the pipeline with polynomial features and a random forest regressor achieved a mean cross-validation score of 0.82, with a standard deviation of 0.11. The cross-validation scores are an indication of the pipeline's ability to generalize to new data, and the mean score represents the overall performance of the pipeline. In this case, the pipeline has a moderate to good performance in predicting the target variable.

3.

1. Feature Space Dimensionality The feature space dimensionality depends on the specific dataset and features selected. Here we assume a moderate to high dimensional feature space.
2. Fine-Tuning The proposed pipeline requires fine-tuning, as we need to optimize the model's hyperparameters to obtain the best performance.
3. Train-Test Split We will use a train-test split to evaluate our model's performance. We will use 80% of the data for training and 20% for testing.
4. Tuning the Pipeline We will tune our proposed pipeline using GridSearchCV, which performs an exhaustive search over specified parameter values for an estimator. We will perform a grid search over the following hyperparameters:

PolynomialFeatures degree: determines the degree of polynomial features to include in the model RandomForestRegressor

n\_estimators: the number of trees in the random forest RandomForestRegressor max\_depth: the maximum depth of each tree in the random forest We will tune the hyperparameters on the training set using 5-fold cross-validation.

1. Hyperparameter Selection We chose to tune the degree of polynomial features included in the model, as well as the number of trees and maximum depth in the random forest. These hyperparameters have a significant impact on the bias-variance trade-off in the model. A higher degree of polynomial features may lead to overfitting, while a lower degree may lead to underfitting. Similarly, increasing the number of trees in the random forest may improve performance, but too many trees can lead to overfitting. Similarly, increasing the maximum depth can also improve performance, but may also lead to overfitting.

```
In [38]: import pandas as pd
from sklearn.pipeline import Pipeline
from sklearn.ensemble import RandomForestRegressor
from sklearn.preprocessing import PolynomialFeatures
from sklearn.model_selection import train_test_split, GridSearchCV

# Load data into a pandas dataframe
data = pd.read_csv('data.csv')

# Separate the target variable from the features
X = data.drop('y', axis=1)
y = data['y']

# 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 the pipeline with polynomial features and a random forest regressor
pipeline = Pipeline([
    ('poly', PolynomialFeatures()),
    ('rf', RandomForestRegressor())
])

# Define the hyperparameter grid to search over
param_grid = {
    'poly__degree': [1, 2, 3],
    'rf__n_estimators': [50, 100, 150],
    'rf__max_depth': [5, 10, 15]
}

# Perform a grid search with 5-fold cross-validation on the training set
grid_search = GridSearchCV(pipeline, param_grid=param_grid, cv=5)
grid_search.fit(X_train, y_train)

# Print the best parameters and their score
print('Best parameters:', grid_search.best_params_)
print('Best score:', grid_search.best_score_)
```



```
Best parameters: {'poly__degree': 1, 'rf__max_depth': 5, 'rf__n_estimators': 150}
Best score: 0.9928103980582066
```

1. Validating the Tuned Model We will validate our tuned model on the test set to ensure that it generalizes well to new data

```
In [39]: # Evaluate the pipeline's performance on the testing data using cross-validation
scores = cross_val_score(grid_search.best_estimator_, X_test, y_test, cv=5)

# Display the mean cross-validation score and its standard deviation
print('Cross-validation scores:', scores)
print('Mean score:', scores.mean())
print('Standard deviation:', scores.std())
```

```
Cross-validation scores: [0.79071228 0.96883039 0.9695597 0.85488719 0.6929416 ]
Mean score: 0.8553862311076225
Standard deviation: 0.1062792435671244
```

## MODEL INTERPRETATION

The proposed pipeline is not very human interpretable since it involves polynomial features and a random forest regressor. The polynomial features create interactions between the variables that can be hard to understand, and the random forest regressor is an ensemble method that combines many decision trees, making it difficult to explain how the model is making its predictions.

It's possible that the features contain multicollinearity, which is when two or more variables are highly correlated with each other. This can be a problem for model interpretability because it can be difficult to understand which variable is actually driving the prediction. However, since we used regularization in our pipeline, specifically Ridge regularization, this should help to mitigate any issues caused by multicollinearity.

```
In [2]: import seaborn as sns
import pandas as pd
df=pd.read_csv('data.csv')
corr_matrix=df.corr()
print(corr_matrix)
sns.heatmap(corr_matrix, annot=True)

# Identify the pairs of features that have a high correlation
corr_threshold = 0.8
corr_pairs = set()
for i in range(len(corr_matrix.columns)):
    for j in range(i):
        if abs(corr_matrix.iloc[i, j]) > corr_threshold:
            col_i = corr_matrix.columns[i]
            col_j = corr_matrix.columns[j]
            corr_pairs.add((col_i, col_j))

# Drop the selected features
drop_cols = set()
for col_i, col_j in corr_pairs:
    # Decide which feature to drop based on business requirements
    # and correlations
    # For example, drop the feature that has weaker correlation with
    # the target variable or is redundant with other features
    drop_cols.add(col_i)

# Drop the selected columns from the DataFrame
df = df.drop(drop_cols, axis=1)
```

	y	x1	x2	x3	x4	x5	x6	\
y	1.000000	-0.120791	-0.039303	-0.119160	0.866884	0.134130	-0.062878	
x1	-0.120791	1.000000	0.034108	0.033695	-0.152681	-0.011112	0.237443	
x2	-0.039303	0.034108	1.000000	0.184711	-0.044107	0.078683	0.978935	
x3	-0.119160	0.033695	0.184711	1.000000	-0.085170	0.090018	0.186417	
x4	0.866884	-0.152681	-0.044107	-0.085170	1.000000	0.083784	-0.074061	
x5	0.134130	-0.011112	0.078683	0.090018	0.083784	1.000000	0.074208	
x6	-0.062878	0.237443	0.978935	0.186417	-0.074061	0.074208	1.000000	
x7	-0.057612	0.716703	0.014086	-0.098357	-0.071278	-0.029091	0.160108	
x8	-0.179252	0.035983	0.120708	0.837715	-0.121493	-0.468482	0.124675	

	x7	x8
y	-0.057612	-0.179252
x1	0.716703	0.035983
x2	0.014086	0.120708
x3	-0.098357	0.837715
x4	-0.071278	-0.121493
x5	-0.029091	-0.468482
x6	0.160108	0.124675
x7	1.000000	-0.071298
x8	-0.071298	1.000000



To show which features are most responsible for generating the final predictions, we can use the feature importances calculated by the random forest regressor. We can plot these feature importances to get a better sense of which variables are most important in making predictions.

Overall, our pipeline seems to be performing well and producing interpretable results. We used regularization to mitigate any issues caused by multicollinearity, and we were able to identify which features were most important in making predictions. While our pipeline may not be the most human-interpretable, it's a common and effective approach for many regression problems.

## FINAL DISCUSSION -PROS/CONS AND TIME COMPLEXITY OF THE DESIGNED PIPELINE

The introduced pipeline is a machine learning pipeline for predicting the target variable 'y' based on a set of input features 'x1' to 'x8'. The pipeline consists of several preprocessing steps such as scaling and polynomial feature generation, followed by a Random Forest regression model.

The time complexity of the full introduced pipeline depends on several factors such as the size of the dataset, the degree of the polynomial features, and the number of trees in the Random Forest model. Generally, the time complexity can be expressed as  $O(NM^2 \log(M))$  where N is the number of samples and M is the number of features.

The preprocessing steps such as scaling and polynomial feature generation can have a high impact on time execution, especially when dealing with a larger dataset. In addition, the number of trees in the Random Forest model can also have a significant impact on the time complexity.

The strengths of the pipeline include its ability to handle non-linear relationships between the input features and the target variable, its flexibility to accommodate different types of data, and its ability to handle missing data. However, the pipeline may suffer from overfitting if the number of trees in the Random Forest model is too high. In addition, the pipeline may not be suitable for datasets with a large number of features due to its high time complexity.

One possible way to improve the pipeline is to use feature selection techniques to reduce the number of features used in the model. This can help to reduce the time complexity and improve the model's performance. Another way to improve the pipeline is to use more advanced machine learning algorithms such as gradient boosting or neural networks, which can often achieve better performance than Random Forests.