

# Section 1

## 1.1 Describe and load your dataset

### Dataset: HR Analytics Employee Attrition & Performance

#### Description:

The dataset contains information about Atlas Lab employees. I'll be working with an employee dataset, which comprises information about employees with the goal of analysing attrition, performance and other aspects of employee retention and workplace dynamics. This type is used in HR analytics to find out traits that define performance and resulting engagement in the workplace. Details such as job, salary, distance from work and tenure encompass demographics of the employees which make the data set useful in analyzing trends in attrition and areas of enhancement in HR practices.

Link: <https://www.kaggle.com/datasets/mahmoudemadabdallah/hr-analytics-employee-attrition-and-performance>

## 1.2. Explain the features:

- EmployeeID:** The recorded date is Unique identifier for each employee.
- FirstName:** The first name of the employee.
- LastName:** The last name of the employee.
- Gender:** The gender of the employee.
- Age:** The age of the employee.
- BusinessTravel:** The frequency of business travel for the employee.
- Department:** The department in which the employee works.
- DistanceFromHome(KM):** The distance between the employee's home and workplace in kilometers.
- State:** The state in which the employee resides.
- Ethnicity:** The ethnicity of the employee.
- MaritalStatus:** The marital status of the employee.
- Salary:** The annual salary of the employee.
- StockOptionLevel:** The level of stock options granted to the employee.
- OverTime:** Whether the employee works overtime (Yes/No).
- HireDate:** The date the employee was hired.
- Attrition:** Whether the employee has left the company (Yes/No).
- YearsAtCompany:** The number of years the employee has been with the company.
- YearsInMostRecentRole:** The number of years the employee has been in their most recent role.
- YearsSinceLastPromotion:** The number of years since the employee's last promotion.
- YearsWithCurrManager:** The number of years the employee has worked with their current manager.

```
In [1]: # Importing libraries
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from scipy import stats
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LinearRegression, Ridge, Lasso, LassoCV
from sklearn.metrics import mean_squared_error
from sklearn.model_selection import train_test_split, KFold, cross_val_score
from scipy.stats import tscore
import statsmodels.api as sm
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import PolynomialFeatures
from sklearn.pipeline import make_pipeline
from scipy.interpolate import LSQUnivariateSpline

# 1.3. Loading the dataset
df = pd.read_csv('Employee.csv')
```

```
In [2]: # 1.4. Show the loaded features
df.head()
```

```
Out[2]:
```

	EmployeeID	FirstName	LastName	Gender	Age	BusinessTravel	Department	DistanceFromHome (KM)	State	Ethnicity	...	MaritalStatus
0	3012-1A41	Leonelle	Simco	Female	30	Some Travel	Sales	27	IL	White	...	Divorced
1	CBCB-9C9D	Leonerd	Aland	Male	38	Some Travel	Sales	23	CA	White	...	Sing
2	95D7-1CE9	Ahmed	Sykes	Male	43	Some Travel	Human Resources	29	CA	Asian or Asian American	...	Married
3	47A0-559B	Ermentrude	Berrie	Non-Binary	39	Some Travel	Technology	12	IL	White	...	Married
4	42CC-040A	Stace	Savege	Female	29	Some Travel	Human Resources	29	CA	White	...	Sing

5 rows x 23 columns

```
In [3]: # Show the loaded features
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1470 entries, 0 to 1469
Data columns (total 23 columns):
#   Column                Non-Null Count  Dtype
---  --
0   EmployeeID             1470 non-null   object
1   FirstName              1470 non-null   object
2   LastName               1470 non-null   object
3   Gender                 1470 non-null   object
4   Age                   1470 non-null   int64
5   BusinessTravel         1470 non-null   object
6   Department             1470 non-null   object
7   DistanceFromHome (KM)  1470 non-null   int64
8   State                  1470 non-null   object
9   Ethnicity              1470 non-null   object
10  Education              1470 non-null   int64
11  EducationField          1470 non-null   object
12  JobRole                1470 non-null   object
13  MaritalStatus          1470 non-null   object
14  Salary                 1470 non-null   int64
15  StockOptionLevel       1470 non-null   int64
16  OverTime               1470 non-null   object
17  HireDate               1470 non-null   object
18  Attrition              1470 non-null   object
19  YearsAtCompany         1470 non-null   int64
20  YearsInMostRecentRole  1470 non-null   int64
21  YearsSinceLastPromotion 1470 non-null   int64
22  YearsWithCurrManager   1470 non-null   int64
dtypes: int64(9), object(14)
memory usage: 244.3+ KB
```

```
In [4]: # 1.5 Show some trend statistics
df.describe()
```

```
Out[4]:
```

	Age	DistanceFromHome (KM)	Education	Salary	StockOptionLevel	YearsAtCompany	YearsInMostRecentRole	YearsSinceLastPromotion	YearsWithCurrManager
count	1470.000000	1470.000000	1470.000000	1470.000000	1470.000000	1470.000000	1470.000000	1470.000000	1470.000000
mean	28.989796	22.502721	2.912925	112956.497959	0.793878	4.562585	2.293197	2.539093	2.293197
std	7.993055	12.811124	1.024165	103342.889222	0.852077	3.288048	2.539093	2.539093	2.539093
min	18.000000	1.000000	1.000000	20387.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	23.000000	12.000000	2.000000	43580.500000	0.000000	2.000000	0.000000	0.000000	0.000000
50%	26.000000	22.000000	3.000000	71999.500000	1.000000	4.000000	1.000000	1.000000	1.000000
75%	34.000000	33.000000	4.000000	142055.750000	1.000000	7.000000	4.000000	4.000000	4.000000
max	51.000000	45.000000	5.000000	547204.000000	3.000000	10.000000	10.000000	10.000000	10.000000

## Section 2

### Lasso Regularization

```
In [5]: # Section 2
# 2.1. Prepare data for Linear regression and lasso regularization.
# Filter the list of features
features2 = ['Age', 'DistanceFromHome (KM)', 'Salary', 'StockOptionLevel', 'YearsAtCompany', 'YearsInMostRecentRole', 'YearsSinceLastPromotion', 'YearsWithCurrManager']

# Define features and target
X = df[features2]
y = df['Salary']

# Standardize the data
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# 2.3. Train linear regression
linear_reg = LinearRegression()
linear_reg.fit(X_scaled, y)

# 2.4. Build a lasso regularization model and choose a value for alpha hyper-parameter.
alpha = 100
lasso_reg = Lasso(alpha=alpha)
lasso_reg.fit(X_scaled, y)

lasso_cv = LassoCV(cv=5, random_state=0) # 5-fold cross-validation
lasso_cv.fit(X_scaled, y)
alpha = lasso_cv.alpha_
print(f"Optimal alpha for Lasso: {alpha}")

# Calculate MSE for both models
y_pred_linear = linear_reg.predict(X_scaled)
y_pred_lasso = lasso_reg.predict(X_scaled)

mse_linear = mean_squared_error(y, y_pred_linear)
mse_lasso = mean_squared_error(y, y_pred_lasso)

print(f"Mean Squared Error for Linear Regression: {mse_linear}")
print(f"Mean Squared Error for Lasso Regression (alpha={alpha}): {mse_lasso}")

# 2.5. Obtain the coefficients for both models
print("\nLinear Regression Coefficients (Scaled):")
for feature, coef in zip(features2, linear_reg.coef_):
    print(f"{feature}: {coef}")

print("\nLasso Regression Coefficients (Scaled):")
for feature, coef in zip(features2, lasso_reg.coef_):
    print(f"{feature}: {coef}")
```

Optimal alpha for Lasso: 103.30773259953392  
Mean Squared Error for Linear Regression: 8.944876800177779e-21  
Mean Squared Error for Lasso Regression (alpha=103.30773259953392): 10000.00000000024

Linear Regression Coefficients (Scaled):  
Age: -1.00683526345181e-10  
DistanceFromHome (KM): -1.4551915228366852e-11  
Salary: 103307.7325995339  
StockOptionLevel: 1.2505552149377763e-12  
YearsAtCompany: 2.546585164964199e-11  
YearsInMostRecentRole: 2.9103930456733704e-11  
YearsSinceLastPromotion: 3.637978807091713e-11  
YearsWithCurrManager: 1.8189894035458565e-11

Lasso Regression Coefficients (Scaled):  
Age: 0.0  
DistanceFromHome (KM): 0.0  
Salary: 103207.73259953388  
StockOptionLevel: 0.0  
YearsAtCompany: 0.0  
YearsInMostRecentRole: 0.0  
YearsSinceLastPromotion: 0.0  
YearsWithCurrManager: 0.0

## 2.6. Analyze the differences between both models.

The coefficients in the linear regression model show extremely small values for most features, but for salary, it has a large value. It means for this model, salary is very important where other features really don't matter much. On the other hand, for lasso regression, only salary has a non-zero coefficient where other features are zero. Lasso regression also thinks the same that salary is important and completely ignores other features.

For the feature selection, linear regression uses all the features without any penalty, so it keeps everything even if some features don't help that much. Lasso regression only picks important features by making the coefficient zero.

The mean squared error for linear regression is almost zero, which means it fits the training data perfectly. However, this can be a sign that the model is overfitting, meaning it might not work well on new, unseen data. Lasso regression has higher Mean Squared Error means it doesn't fit the training data as perfectly. But because it ignores unimportant features, it may perform better on new data by avoiding overfitting.

Overall, linear regression offers highly accurate fit on training, but it can be overfitting because it includes all features, even the unimportant ones. Lasso regression simplifies the model by removing features that don't add much value. It might not fit the training data as perfectly, but it's more likely to work better with new data because it avoids overfitting.

## Section 4

### K-fold Cross-validation

```
In [6]: # Section 4 - K-fold Cross-validation

# 4.1. Set up your data for K-fold cross-validation
features2 = ['Age', 'DistanceFromHome (KM)', 'Education', 'StockOptionLevel', 'YearsAtCompany', 'YearsInMostRecentRole', 'YearsSinceLastPromotion', 'YearsWithCurrManager']
X = df[features2]
y = df['Salary'] # MAKE MODIFICATIONS

# 4.2. Initialize the Linear Regression model
linear_reg = LinearRegression()

# 4.3. Define a 5-fold cross-validation split
kf = KFold(n_splits=5, shuffle=True, random_state=42) # MAKE MODIFICATIONS k = 5 or other

# 4.5. Manually loop through each fold
fold = 1
mse_scores = [] # List to store the MSE for each fold
for train_index, test_index in kf.split(X):
    # Split data into training and testing sets for the current fold
    X_train, X_test = X.iloc[train_index], X.iloc[test_index]
    y_train, y_test = y.iloc[train_index], y.iloc[test_index]

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

    # Predict on the testing set
    y_pred = linear_reg.predict(X_test)

    # Calculate MSE for the current fold
    mse = mean_squared_error(y_test, y_pred)
    mse_scores.append(mse)

print(f"Fold {fold}: Mean Squared Error = {mse}")
fold += 1

# Print results
average_mse = np.mean(mse_scores)
variance_mse = np.var(mse_scores) # Calculate variance of MSE across folds

print("\nCross-validation results:")
print(f"Mean MSE across all folds: {average_mse}")
print(f"Variance of MSE across all folds: {variance_mse}")
print(f"MSE scores for each fold: {mse_scores}")

# Normalize MSE scores by dividing by the maximum MSE (This makes it easy to interpret)
max_mse = max(mse_scores)
normalized_mse_scores = [mse / max_mse for mse in mse_scores]
normalized_mean_mse = average_mse / max_mse
normalized_variance_mse = variance_mse / (max_mse ** 2)

# Display normalized results
print("\nCross-validation results (Normalized):")
for i, mse in enumerate(normalized_mse_scores, 1):
    print(f"Fold {i}: Normalized MSE = {mse:.4f}")

print(f"Normalized Mean MSE across all folds: {normalized_mean_mse:.4f}")
print(f"Normalized Variance of MSE across all folds: {normalized_variance_mse:.4f}")
print(f"Normalized MSE scores for each fold: {normalized_mse_scores}")
```

Fold 1: Mean Squared Error = 9550137466.174374  
Fold 2: Mean Squared Error = 6615679659.494051  
Fold 3: Mean Squared Error = 6123104017.565921  
Fold 4: Mean Squared Error = 6559927084.609264  
Fold 5: Mean Squared Error = 7178058579.289058

Cross-validation results:  
Mean MSE across all folds: 7205381361.426534  
Variance of MSE across all folds: 1.4868622589867497e+18  
MSE scores for each fold: [9550137466.174374, 6615679659.494051, 6123104017.565921, 6559927084.609264, 7178058579.289058]

Cross-validation results (Normalized):  
Fold 1: Normalized MSE = 1.0000  
Fold 2: Normalized MSE = 0.6927  
Fold 3: Normalized MSE = 0.6412  
Fold 4: Normalized MSE = 0.6869  
Fold 5: Normalized MSE = 0.7516

Normalized Mean MSE across all folds: 0.7545  
Normalized Variance of MSE across all folds: 0.0163  
Normalized MSE scores for each fold: [1.0, 0.692731341743376, 0.641153495355783, 0.6868934722504118, 0.7516183515381867]

## 4.6. Analysis:

The cross validation outcomes revealed that the model's MSE is different with the five folds, where fold 1 has the highest error rate of 0.1647 and fold 3 has the lowest of 0.1289. It seems that the libraries chosen are close to being final, so the average MSE is in the region of 7.2 billion, and a high variance up to 5 billion indicates instability in the choice of data splits. Similar to equation 4, the normalized MSE ranges the Fold 1 at a level of 1 indicating that it stands in for the actual meaning of a poor model while other folds that show comparatively lower errors depict comparatively better performance. It remains at around 0.75 for average normalized MSE, while the variance is 0.0163, implying that there is some deviation, but not highly contrasting results among the folds. Overall, it is a good model for some of the comparisons but has shown fairly mixed results as seen in the table.

```
In [7]: # Section 5 - Linear Model Selection Statistics

# 5.1. Set up your data for linear regression # MAKE MODIFICATIONS
features2 = ['Age', 'DistanceFromHome (KM)', 'StockOptionLevel', 'YearsAtCompany', 'YearsInMostRecentRole', 'YearsSinceLastPromotion', 'YearsWithCurrManager']
X = df[features2]
y = df['Salary']

# 5.2. Define function to compute Cp
def compute_cp(model, X, y):
    mse = np.mean((model.predict(X) - y) ** 2)
    p = len(model.params) - 1 # Number of predictors
    n = len(y)
    cp = mse * 2 * p * mse / (n - p - 1)
    return cp

# 5.3. Compute metrics for models with increasing numbers of predictors
predictors = X.columns
cp_values, aic_values, bic_values, adjr2_values = [], [], [], []

for k in range(1, len(predictors) + 1):
    chosen_predictors = predictors[:k]
    X_subset = X[chosen_predictors]
    X_subset = sm.add_constant(X_subset) # Add constant for intercept
    model = sm.OLS(y, X_subset).fit()

    cp_values.append(compute_cp(model, X_subset, y))
    aic_values.append(model.aic)
    bic_values.append(model.bic)
    adjr2_values.append(model.rsquared_adj)

# 5.4. Plotting
fig, ax1 = plt.subplots(figsize=(12, 6))

ax2 = ax1.twinx()
ax3 = ax1.twinx()
ax4 = ax1.twinx()

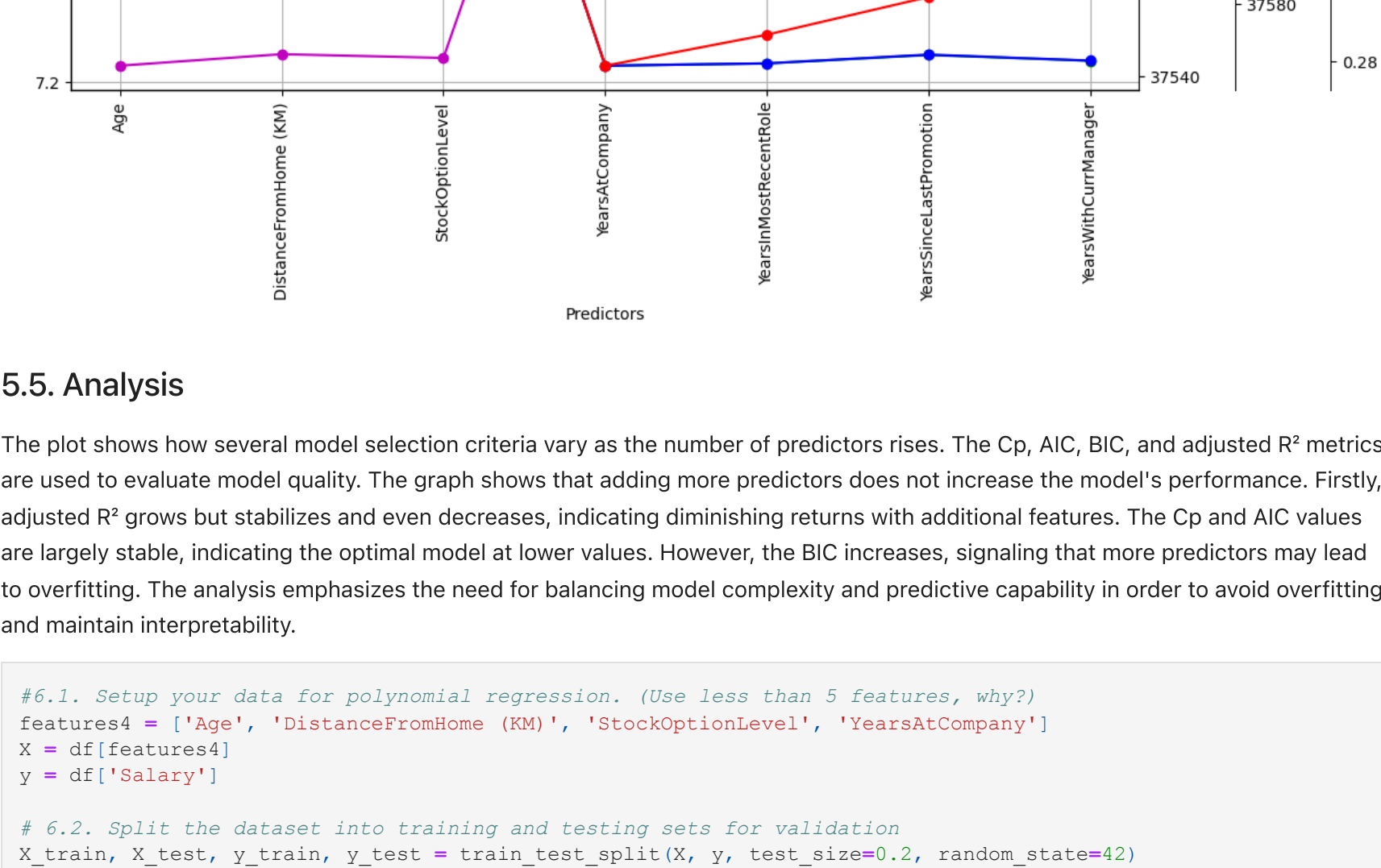
# Adjust the position of the third and fourth axes to avoid overlap
ax3.spines['right'].set_position(('outward', 60))
ax4.spines['right'].set_position(('outward', 120))

# Plot metrics on each appropriate axis
ax1.plot(predictors, cp_values, 'g-', label='Cp', marker='o')
ax2.plot(predictors, aic_values, 'b-', label='AIC', marker='o')
ax3.plot(predictors, bic_values, 'r-', label='BIC', marker='o')
ax4.plot(predictors, adjr2_values, 'm-', label='Adj R^2', marker='o')

# Set labels and title
ax1.set_xlabel('Predictors')
ax1.set_ylabel('Cp', color='g')
ax2.set_ylabel('AIC', color='b')
ax3.set_ylabel('BIC', color='r')
ax4.set_ylabel('Adjusted R^2', color='m')
plt.title('Model Selection Criteria by Number of Predictors')

# Add grid and rotate x-axis labels
ax1.grid(True)
ax1.set_xticks(range(len(predictors)+1))
ax1.set_xticklabels(predictors, rotation=90) # Rotate x-axis labels vertically

# Display
plt.show()
```



## 5.5. Analysis

The plot shows how several model selection criteria vary as the number of predictors rises. The Cp, AIC, BIC, and adjusted R<sup>2</sup> metrics are used to evaluate model quality. The graph shows that adding more predictors does not increase the model's performance. Firstly, adjusted R<sup>2</sup> grows but stabilizes and even decreases, indicating diminishing returns with additional features. The Cp and AIC values are largely stable, indicating the optimal model at lower values. However, the BIC increases, signaling that more predictors may lead to overfitting. The analysis emphasizes the need for balancing model complexity and predictive capability in order to avoid overfitting and maintain interpretability.

```
In [8]: #6.1. Setup your data for polynomial regression. (Use less than 5 features, why?)
features4 = ['Age', 'DistanceFromHome (KM)', 'StockOptionLevel', 'YearsAtCompany']
X = df[features4]
y = df['Salary']

# 6.2. Split the dataset into training and testing sets for validation
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# 6.3. 10th-degree polynomial regression
poly10_model = make_pipeline(PolynomialFeatures(10), LinearRegression())
poly10_model.fit(X_train, y_train)
y_pred_poly10 = poly10_model.predict(X_test)
mse_poly10 = mean_squared_error(y_test, y_pred_poly10)

# 6.4. 4th-degree polynomial regression
poly4_model = make_pipeline(PolynomialFeatures(4), LinearRegression())
poly4_model.fit(X_train, y_train)
y_pred_poly4 = poly4_model.predict(X_test)
mse_poly4 = mean_squared_error(y_test, y_pred_poly4)

# 6.5. 3rd-degree polynomial regression
poly3_model = make_pipeline(PolynomialFeatures(3), LinearRegression())
poly3_model.fit(X_train, y_train)
y_pred_poly3 = poly3_model.predict(X_test)
mse_poly3 = mean_squared_error(y_test, y_pred_poly3)

# 6.6. Display MSE for different polynomial degrees
print(f"Mean Squared Error for 10th-degree Polynomial: {mse_poly10:.2f}")
print(f"Mean Squared Error for 4th-degree Polynomial: {mse_poly4:.2f}")
print(f"Mean Squared Error for 3rd-degree Polynomial: {mse_poly3:.2f}")

# Normalize MSE scores by dividing by the maximum MSE (makes it easier to interpret)
mse_scores = [mse_poly10, mse_poly4, mse_poly3]
average_mse = np.mean(mse_scores)
variance_mse = np.var(mse_scores)
max_mse = max(mse_scores)
normalized_mse_scores = [mse / max_mse for mse in mse_scores]
normalized_mean_mse = average_mse / max_mse
normalized_variance_mse = variance_mse / (max_mse ** 2)

# Display normalized results
print("\nCross-validation results (Normalized):")
for i, mse in enumerate(normalized_mse_scores, 1):
    print(f"Fold {i}: Normalized MSE = {mse:.4f}")

print(f"Normalized Mean MSE across all models: {normalized_mean_mse:.4f}")
print(f"Normalized Variance of MSE across all models: {normalized_variance_mse:.4f}")
print(f"Normalized MSE scores for each model: {normalized_mse_scores}")
```

Mean Squared Error for 10th-degree Polynomial: 13751906782570242.00  
Mean Squared Error for 4th-degree Polynomial: 8241278442.32  
Mean Squared Error for 3rd-degree Polynomial: 8185296431.11

Cross-validation results (Normalized):  
Normalized MSE = 1.0000  
Normalized MSE = 0.0000  
Normalized MSE = 0.0000  
Normalized Mean MSE across all models: 0.3333  
Normalized Variance of MSE across all models: 0.2222  
Normalized MSE scores for each model: [1.0, 0.5992825993242384e-07, 0.592117448529433e-07]

## 6.6. Analysis

The performance of polynomial regression models shows that the 10th degree model is better than the 4th degree and 3rd degree in terms of mean squared error (MSE). This is because the 10th degree polynomial was overfitted to the data and that automatically guided it towards poor performance on the test based set. The reduced MSE and similar performance of the 4th degree and 3rd degree models indicate a better generalization with the use of fewer parameters, as well as stability regarding hyper-parameter tuning. In normalized form, the 10th degree model was taken as a baseline (1.0) and all of the other models had an MSE that normalized to zero indicating their far greater performance compared to the 10th degree model. The mean of 0.3333 shows the normalized MSE across the models, and we see from the variance of 0.2222 that there is a great deal of difference between model performances. In general, less complex polynomial models (3rd and 4th degrees) were found to be more accurate and less sensitive to random variations in data, than a highly complex model 10th degree, even though the complexities were equipped through an automatic search process applied beforehand.



