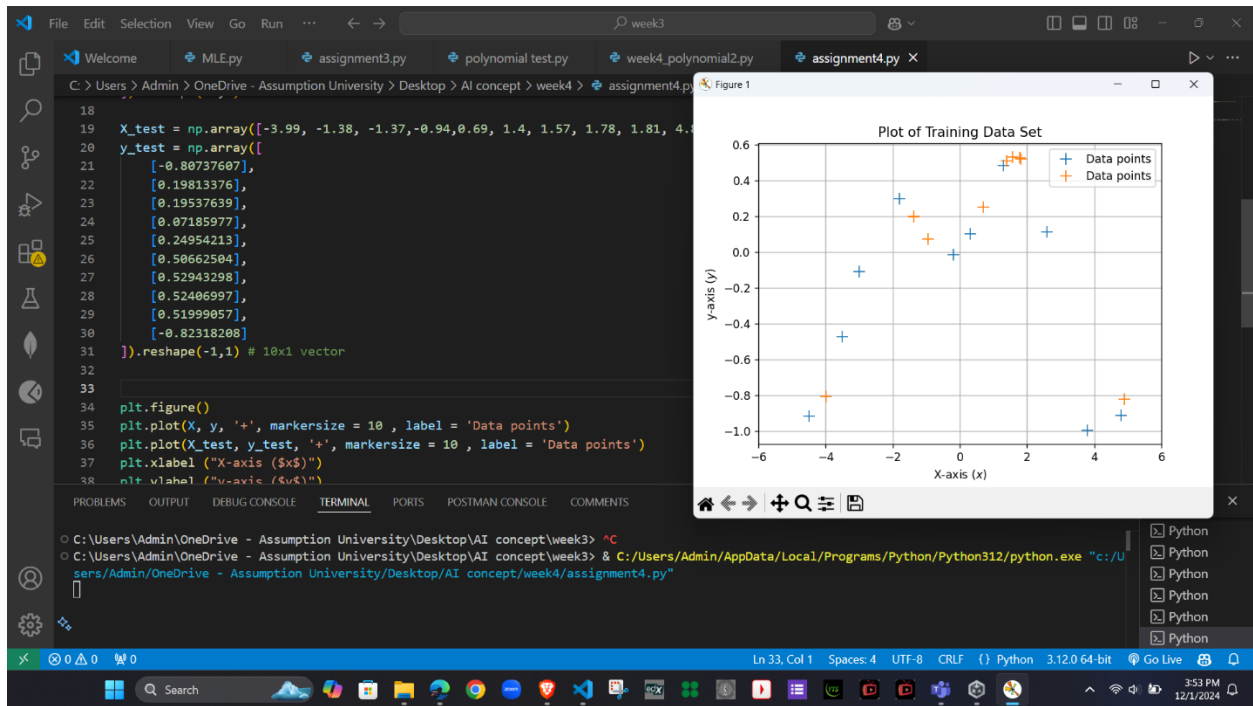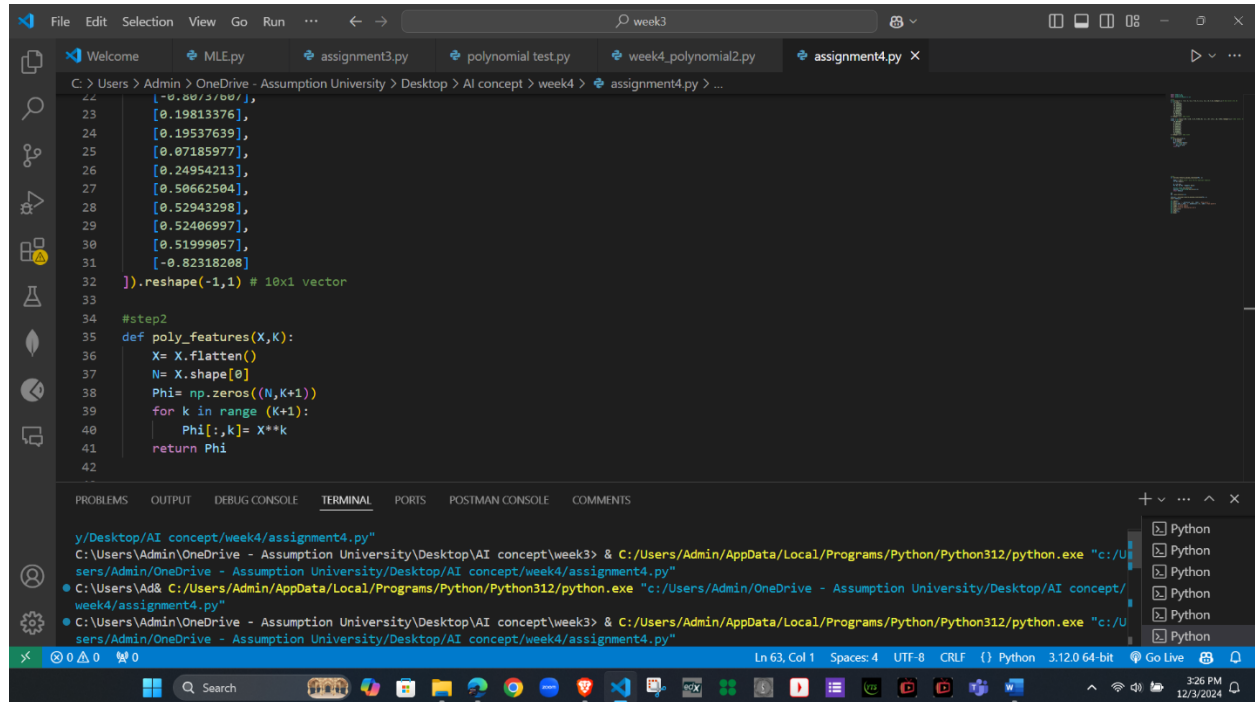# Assignment 4

# Q1. Maximum Likelihood Estimation with Features

# Step 1: Understanding and Plotting the Data
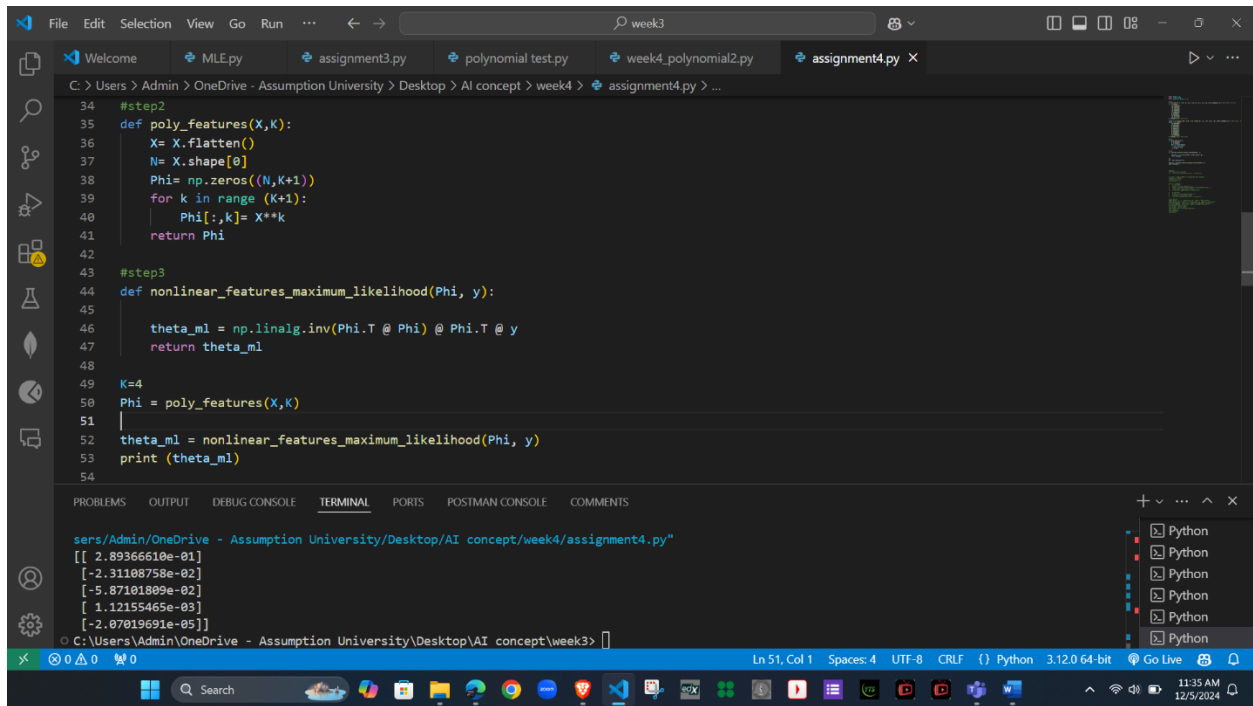
# Step 2: Polynomial Feature Transformation



## The purpose and the effect of this transformation:

Polynomial transformation enables linear models to fit nonlinear data by adding higher-degree features. It improves flexibility but risks underfitting with low degrees or overfitting with high degrees. The right degree balances accuracy and complexity for better generalization.

# Step 3: Fitting the Model Using Maximum Likelihood

# Step 4: Model Evaluation Using RMSE

# Step 5: Selecting the Best Model



Why was this degree optimal?

- The optimal degree minimizes the RMSE for the test dataset, striking a balance between underfitting (too simple) and overfitting (too complex).

- Degrees that are too low (e.g., 0 or 1) underfit, being too simple and cannot catch the small details in the data, so they don't fit well.

- Degrees that are too high (e.g., >10) overfit, capturing noise rather than the true relationship.

Model Complexity and Overfitting/Underfitting:

- Lower degrees: Underfitting, as the model cannot capture complex patterns.

- Higher degrees: Overfitting, as the model tries to fit every data point, including noise.

- Optimal degree: Provides the best generalization for unseen data (test set).

# Step 6: Conclusion

Summary of Findings:

- The polynomial regression model's performance depends heavily on the degree of the polynomial.

- Low-degree models (e.g., 0 or 1): These models were too simple and underfit the data, missing important patterns and resulting in high training and test RMSE.

- High-degree models (e.g., >10): These models overfit the training data, capturing noise instead of meaningful trends, which led to high test RMSE.

- The optimal degree was the one that minimized test RMSE, achieving a balance between underfitting and overfitting.


Trade-offs Between Model Complexity and Prediction Accuracy:

- Low Complexity (Low Degree):

  - Advantages: Simple, easier to interpret, and less likely to overfit.

  - Disadvantages: Poor performance if the relationship between variables is complex (underfitting).

- High Complexity (High Degree):

  - Advantages: Can model complex relationships in the data and fit the training data very well.

- Disadvantages: Risk of overfitting, where the model captures noise instead of general patterns, leading to poor test performance.

- Balanced Complexity (Optimal Degree):

  - The optimal degree balances the trade-off, capturing the true patterns in the data without fitting the noise, resulting in better generalization to unseen data.

# Q2. K-means

# Step 1: Prepare Dataset

- Clean the unnecessary data, duplicate data and missing value

# Step 2: Visualize Data Before Clustering

# Step 3: Implement K-means Clustering and Determine Optimal Clusters

- Code for step3 using Elbow Method Code



```python
#step3
import numpy as np
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.datasets import make_blobs

X, _ = make_blobs(n_samples=300, centers=4, cluster_std=0.60, random_state=0)

wcss = []
for i in range(1, 11):  # Test 1 to 10 clusters
    kmeans = KMeans(n_clusters=i, init='k-means++', max_iter=300, n_init=10, random_state=0)
    kmeans.fit(X)
    wcss.append(kmeans.inertia_)

plt.figure(figsize=(10, 8))
plt.plot(range(1, 11), wcss, marker='o')
plt.title('Elbow Method For Optimal Clusters')
plt.xlabel('Number of Clusters')
plt.ylabel('WCSS (Within-Cluster Sum of Squares)')
plt.grid(True)
plt.show()

optimal_clusters = 5
kmeans = KMeans(n_clusters=optimal_clusters, init='k-means++', max_iter=300, n_init=10, random_state=0)
y_kmeans = kmeans.fit_predict(X)
print(y_kmeans)
```

- Plot the Elbow Method Graph



- K-mean output

# Step 4: Visualize Data After Clustering

# Step 5: Interpret Your Results

The clustering results provide insights into customer segmentation based on Annual Income (k$) and Spending Score (1-100). Here's a detailed breakdown of each cluster:

Cluster Sizes:

- Cluster 0 (81 customers): Largest group with moderate income and spending behavior.

- Cluster 1 (39 customers): High-income, high-spending customers (premium segment).

- Cluster 2 (35 customers): High-income, low-spending customers (potential for engagement).

- Cluster 3 (23 customers): Low-income, low-spending customers (limited purchasing power).

- Cluster 4 (22 customers): Low-income, high-spending customers (price-sensitive but engaged).

Cluster Statistics Analysis:

1. Cluster 0 (Moderate Income, Moderate Spending):

   o Average Annual Income: ~$55.3k

   o Average Spending Score: ~49.5

   o Represents customers with balanced spending habits, likely mid-tier or average-value customers.

   o Opportunity: Maintain their satisfaction with general promotions or loyalty programs.

2. Cluster 1 (High Income, High Spending):

- Average Annual Income: ~$86.5k

- Average Spending Score: ~82.1

- High-value customers who spend a lot and likely purchase premium products/services.

- Opportunity: Strengthen loyalty through exclusive rewards, VIP programs, and personalized experiences.

3. Cluster 2 (High Income, Low Spending):

- Average Annual Income: ~$88.2k

- Average Spending Score: ~17.1

- High-income customers with low engagement or spending behavior.

- Opportunity: Explore why these customers are not spending much (e.g., mismatched offerings, lack of interest). Use targeted marketing to improve engagement.

4. Cluster 3 (Low Income, Low Spending):

- Average Annual Income: ~$26.3k

- Average Spending Score: ~20.9

- Budget-conscious customers with low purchasing power.

- Opportunity: Offer affordable options, discounts, or budget-friendly product lines to attract and retain these customers.

5. Cluster 4 (Low Income, High Spending):

- Average Annual Income: ~$25.7k

- Average Spending Score: ~79.4

- Price-sensitive customers who are highly engaged despite having limited income.

- Opportunity: Provide value-based products and loyalty programs to encourage frequent purchases.

---

How This Information Helps Identify Specific Groups:

1. Targeted Marketing Strategies:

   - Cluster 1: Focus on retaining high-value customers through personalized and exclusive offers.

   - Cluster 2: Design campaigns to convert high-income, low-spending customers into active spenders.

   - Cluster 4: Develop affordable product bundles or discounts to maximize the engagement of price-sensitive customers.

2. Customer Retention:

   - Enhance customer retention in Cluster 0 through loyalty programs that reward consistent spending.

   - Identify disengagement reasons in Cluster 2 and address them with surveys or tailored incentives.

3. Product and Service Design:

   - Offer premium products and services to Cluster 1.

- o Introduce budget-friendly options for Cluster 3 and Cluster 4.

4. Resource Allocation:

- o Allocate more resources to engage and retain high-potential customers in Clusters 1, 2, and 4.

- o Focus less on Cluster 3, as they represent a low-spending group with limited income.

5. Long-Term Strategies:

- o Monitor spending trends in Cluster 2 to identify signs of increased engagement.

- o Encourage upselling and cross-selling to moderate spenders in Cluster 0.

# Additional code