

1. What do you understand about Machine Learning? What are the main stages in the Machine Learning pipeline? Briefly explain the difference between bad data and a bad model.

Machine Learning (ML) is a field of AI where systems learn patterns from data instead of being explicitly programmed. A machine learning model improves its performance automatically as it sees more examples.

Stages in the Machine Learning Pipeline:

1. Data Collection  
Gathering raw data from files, APIs, sensors, etc.
2. Data Cleaning  
Handling missing values, removing duplicates, fixing inconsistencies.
3. Feature Engineering  
Selecting and transforming features so the model can learn better.
4. Model Selection  
Choosing which model to train (regression, decision tree, SVM, etc.).
5. Training the Model  
Feeding data to the model so it learns patterns.
6. Testing & Evaluation  
Checking how well the model performs on unseen data.
7. Deployment  
Sending the model into production (apps, services, systems).

Bad Data	Bad Model
Missing values, noise, wrong labels, outliers	Wrong algorithm choice or poor tuning
Garbage-in-garbage-out	Can't learn patterns even with good data
Leads to inaccurate learning	Leads to poor predictions even with good data

2. Why is a regression model used in Machine Learning? Is it a supervised or unsupervised learning method, and why? What are the models for simple and multiple linear regression? Provide the formulas used to calculate the values of betas (no need to derive them) in both cases.

Regression is used in Machine Learning when you want to predict a continuous value. Examples: house price, temperature, salary, medical test value.

Regression is Supervised Learning because the model is trained using labeled data, which is features which are known outputs.

The formulas used to calculate the values of betas in 3 different types of linear regression are given below:

## Types of Linear Regression

### 1. Simple Linear Regression

○ One independent variable

○ Model:

$$y = \beta_0 + \beta_1 x + \epsilon$$

### Multiple Linear Regression

• Multiple independent variables

• Model:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$$

### Polynomial Regression

• Extends linear regression by including **non-linear powers of X**

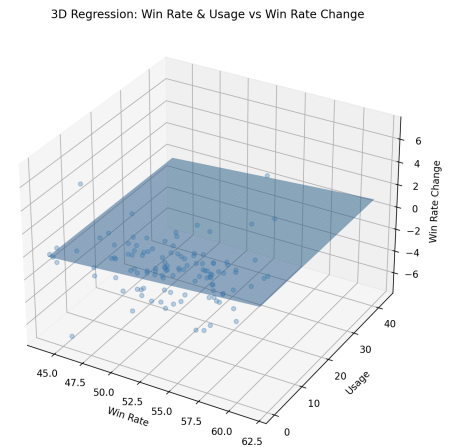
• Equation

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \dots + \beta_n X^n$$

3. Download a real-world dataset (not used in class) to train and test a multiple linear regression model. Use it to predict the output. Clearly mention the dataset you used.

a. First, train the model using only two features. Also, create a scatter plot of the dataset using those two features and overlay the regression model. Include screenshots of your code, output, and scatter plot.

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 from sklearn.linear_model import LinearRegression
4 import numpy as np
5 from mpl_toolkits.mplot3d import Axes3D
6
7 # Load dataset
8 df = pd.read_csv('/Users/ishan-college/Downloads/clash_royale_cards.csv')
9
10 # --- Select TWO features and ONE target ---
11 X = df[["Win Rate", "Usage"]] # features
12 y = df["Win Rate Change"] # target
13 # -----
14
15 # Train model
16 model = LinearRegression()
17 model.fit(X, y)
18
19 print("Intercept ( $\beta_0$ ):", model.intercept_)
20 print("Coefficients ( $\beta$ ):", model.coef_)
21
22 # 3D Scatter Plot
23 fig = plt.figure(figsize=(12, 8))
24 ax = fig.add_subplot(111, projection='3d')
25
26 ax.scatter(X["Win Rate"], X["Usage"], y, alpha=0.3)
27
28 # Create prediction plane
29 x_surf, y_surf = np.meshgrid(
30     np.linspace(X["Win Rate"].min(), X["Win Rate"].max(), 50),
31     np.linspace(X["Usage"].min(), X["Usage"].max(), 50)
32 )
33
34 z_surf = (
35     model.intercept_
36     + model.coef_[0] * x_surf
37     + model.coef_[1] * y_surf
38 )
39
40 ax.plot_surface(x_surf, y_surf, z_surf, alpha=0.5)
41
42 # Labels
43 ax.set_xlabel("Win Rate")
44 ax.set_ylabel("Usage")
45 ax.set_zlabel("Win Rate Change")
46 ax.set_title("3D Regression: Win Rate & Usage vs Win Rate Change")
47
48 plt.show()
49
50 # --- User Prediction ---
51 try:
52     win_rate = float(input("Enter Win Rate: "))
53     usage = float(input("Enter Usage: "))
54
55     input_df = pd.DataFrame([win_rate, usage], columns=["Win Rate", "Usage"])
56     prediction = model.predict(input_df)
57
58     print(f"Predicted Win Rate Change: {prediction[0]:.2f}")
59
60 except ValueError:
61     print("Enter numeric values only.")
62
63
```



```
Intercept ( $\beta_0$ ): -3.3196465317079062
Coefficients ( $\beta$ ): [ 0.07033471 -0.00118427]
Enter Win Rate: 50
Enter Usage: 55
Predicted Win Rate Change: 0.13
```

b. Next, train the model using more than two features of your choice. Show the calculated beta values and use the model to predict the output for some input values.

```
1 import pandas as pd
2 from sklearn.linear_model import LinearRegression
3
4 # Load dataset
5 df = pd.read_csv('/Users/ishan-college/Downloads/clash_royale_cards.csv')
6
7 # ---- Select MORE THAN 2 features ----
8 X = df[["Win Rate", "Usage", "Usage Change"]]
9
10 # Select target variable
11 y = df["Win Rate Change"]
12 # -----
13
14 # Train model
15 model = LinearRegression()
16 model.fit(X, y)
17
18 # Print intercept and coefficients
19 print("Intercept ( $\beta_0$ ):", model.intercept_)
20 print("\nCoefficients ( $\beta$  values):")
21 for feature, coef in zip(X.columns, model.coef_):
22     print(f"{feature}: {coef}")
23
24 # ---- Make Predictions ----
25 print("\n--- Prediction Example ---")
26
27 try:
28     win_rate = float(input("Enter Win Rate: "))
29     usage = float(input("Enter Usage: "))
30     usage_change = float(input("Enter Usage Change: "))
31
32     input_df = pd.DataFrame([win_rate, usage, usage_change],
33                             columns=["Win Rate", "Usage", "Usage Change"])
34
35     prediction = model.predict(input_df)
36
37     print(f"\nPredicted Win Rate Change: {prediction[0]:.4f}")
38
39 except ValueError:
40     print("Invalid input. Please enter numeric values.")
41
```

```
Intercept ( $\beta_0$ ): -2.3304773484490116
Coefficients ( $\beta$  values):
Win Rate: 0.05133393960423274
Usage: 0.0055068847975453476
Usage Change: 0.5892472048789782

--- Prediction Example ---
Enter Win Rate: 5
Enter Usage: 7
Enter Usage Change: 1
Predicted Win Rate Change: -1.4460
```

For this dataset, I used a real-world Clash Royale balance dataset named: clash\_royale\_cards.csv from [Clash Royale Cards Stats - Kaggle](#)