

1. Data Gathering & Description

For this project, we will be using two original datasets to compare neural networks and classical machine learning methods. Let's just introduce the first dataset and test a few things before we get into the second dataset.

Our first dataset, covering soccer/football match statistics in various top-tier and second-tier leagues in various countries, was sourced via the **football-data.org API**. We signed up for a free API token, which was used to authenticate requests, so we could obtain historical match data from various leagues in Europe (such as the Bundesliga and the Premier League) as well as the top league in Brazil.

For this project, to fulfill our needs and requirements, we decided to pull data from two full seasons' worth of data from various leagues.

A Python script was used to pull the data programmatically, including match results, team information, and goal statistics. Here's a sample of the raw JSON output, illustrating the structure of the data we planned to work with.

```
{
  "head2head": {
    "number Of Matches": 1,
    "totalGoals": 4,
    ...
  },
  "match": {
    "id": 200282,
    "competition": {
      "id": 2001,
      "name": "UEFA Champions League"
    },
    ...
  },
  "homeTeam": {
    "id": 86,
    "name": "Real Madrid CF",
    ...
  },
  "awayTeam": {
    "id": 64,
    "name": "Liverpool FC",
    ...
  }
}
```

```
        },
        "score": {
            "winner": "HOME_TEAM",
            "duration": "REGULAR",
            "fullTime": {
                "homeTeam": 3,
                "awayTeam": 1
            },
            ...
        }
    }
}
```

Engineering Challenges:

With the help of Gemini AI and ChatGPT, the newer structured utilized by football-data.org didn't create too much of barrier in navigating the API's version changes. We copied and pasted the documentation, along with the request procedure, and we had plenty of code to work on our project.

Gemini helped us: *parse and synthesize* the structure and content of the documentation for each resource such as Match, Team, or Person; *identify relevant JSON fields*, such as score and homeTeam within the Match resource, that would be the most useful for our regression problem; and *formulate a data strategy* for iterating through seasons and competitions to collect a large enough dataset.

The biggest concern was managing the API's rate limits (for a free subscriber) by using a somewhat delayed loop. This way, we could pull the data for multiple seasons - to meet the 3,000-example requirement - while still making our requests in a time-efficient manner. This also helps us stretch the requests out enough, so we wouldn't make an excessive number of requests, resulting in a lovely HTTP429 ("too many requests") error.

A Sports Betting Analytics Test

The soccer API data, as we started to look through it, could be a simplified version of a **sports prediction model**. So what if our dataset was initially used as a betting analytics model - for this part of the project?

Here's why:

1. **Historical Data:** Patterns are usually the key! Both our basic model and a professional analytics model relies on a large volume of historical data. The API's ~~detailed records of past matches, scores, and team stats are the core of this~~

Detailed records of past matches, scores, and team stats are the core of this similarity.

2. **Feature Engineering:** One of uses a betting site, and in the football/soccer section, there's a section in upcoming matches that discusses any past head-to-head matches and the results of the last 5 games (for each team). That information is a fancy *feature that's was engineered* - and we can do that, too.
3. **Model Building:** We will train a model (such as a **Logistic Regression** or **Random Forest**) on these features to predict an outcome. The model (we hope) will weigh each feature's importance, such as how a team's home advantage is a more significant predictor than its average goals scored.
4. **Prediction:** Once trained, our model can take the features for an upcoming match and hopefully produce a fairly accurate prediction, which is often a probability (e.g., "Team A has a 65% chance of winning"). This is the key output for a basic betting analysis.

Although this part is a basic model for our project, the methodology used will likely follow what plenty of analysts use, including some of strategies and thought processing that Vegas uses to set their odds.

▼ The Model and its Features

We identified several features to focus on for our **classical machine learning model**. They cover crucial aspects of a soccer match:

- **Home/Away & Against Which Team:** These are fundamental. Home-field advantage and the opponent's strength are two of the most important predictors of a match's outcome.
- **Goalie, Yellow/Red Cards:** These are key tactical and performance features, since the decision to use a particular goalie or the aftermath of a red card can significantly affect a match. Our model would learn how these factors affect the final score.
- **Time of Year/Month/Day/Night:** We thought this is something that may alter the match, such as changes in weather or daylight vs stadium lighting. This could also highlight player fatigue at the end of a long season as well as determining if certain teams have a scheduling advantage.

What the Model Can Predict

We expect our model to do exactly this: **predict head-to-head match outcomes** during the season. This would be a **classification problem** where the model's output is one of

the season. This would be a **Classification problem** where the model's output is one of three categories: Home Team Wins, Away Team Wins, or Draw.

By using home/away status, goalie performance, and yellow/red card data as features, we're looking at creating a model that's based on data among several leagues. That, hopefully, will also assist us with a model that can predict match outcomes across different leagues.

We expect the model to highlight patterns that are common to all of soccer, such as the impact of a red card or home-field advantage, rather than just patterns specific to one league. For example, the Bundesliga in Germany is known to be a bit lenient when it comes to aggressive playing and physical interaction, so the yellow/red card numbers may not affect it like the penalties do in the Premier League in the UK.

```
import requests
import pandas as pd
import time

# YOUR API KEY GOES HERE
api_key = "e0d3c18b3e2c496e9a79bf92ccbbe53f"
BASE_URL = "https://api.football-data.org/v4/"

headers = {
    'X-Auth-Token': api_key,
    'X-Unfold-Goals': 'true'
}

# The leagues to include, based on your free plan
competition_codes = ['BL1', 'DED', 'BSA', 'PD', 'FL1', 'ELC', 'PPL', 'SA', 'PL']
seasons = [2023, 2024]

all_matches = []
unique_match_ids = set()

print("Starting data collection for multiple leagues and seasons...")
for competition_code in competition_codes:
    for season in seasons:
        print(f"\nFetching matches for {competition_code} in {season}/{season+1}")
        url = f"{BASE_URL}competitions/{competition_code}/matches?season={seasc

try:
    response = requests.get(url, headers=headers)
    response.raise_for_status()
    data = response.json()

    if 'matches' in data:
        for match in data['matches']:
            if match['id'] not in unique_match_ids:
```

```
        all_matches.append(match)
        unique_match_ids.add(match['id'])
    print(f" Successfully fetched {len(data['matches'])} matches.")
else:
    print(f" No matches found for {competition_code} in that seasc

except requests.exceptions.RequestException as e:
    print(f" An error occurred: {e}")

time.sleep(15) # Respect the API rate limit

# --- Create a Final DataFrame ---
print("\nCreating a single DataFrame from the collected data...")
if all_matches:
    final_df = pd.DataFrame(all_matches)

    print(f"\nTotal unique matches in the final DataFrame: {len(final_df)}")

    print("\nDataFrame created. Here are the first 5 rows:")
    print(final_df.head())

    print("\nDataFrame Columns:")
    print(final_df.columns)
else:
    print("No match data was collected.")
```

Starting data collection for multiple leagues and seasons...

Fetching matches for BL1 in 2023/2024 season...
Successfully fetched 306 matches.

Fetching matches for BL1 in 2024/2025 season...
Successfully fetched 306 matches.

Fetching matches for DED in 2023/2024 season...
Successfully fetched 306 matches.

Fetching matches for DED in 2024/2025 season...
Successfully fetched 306 matches.

Fetching matches for BSA in 2023/2024 season...
Successfully fetched 380 matches.

Fetching matches for BSA in 2024/2025 season...
Successfully fetched 380 matches.

Fetching matches for PD in 2023/2024 season...
Successfully fetched 380 matches.

Fetching matches for PD in 2024/2025 season...
Successfully fetched 380 matches.

Fetching matches for FL1 in 2023/2024 season...

```
Successfully fetched 306 matches.

Fetching matches for FL1 in 2024/2025 season...
Successfully fetched 306 matches.

Fetching matches for ELC in 2023/2024 season...
Successfully fetched 557 matches.

Fetching matches for ELC in 2024/2025 season...
Successfully fetched 557 matches.

Fetching matches for PPL in 2023/2024 season...
Successfully fetched 306 matches.

Fetching matches for PPL in 2024/2025 season...
Successfully fetched 306 matches.

Fetching matches for SA in 2023/2024 season...
Successfully fetched 380 matches.

Fetching matches for SA in 2024/2025 season...
Successfully fetched 380 matches.

Fetching matches for PL in 2023/2024 season...
Successfully fetched 380 matches.

Fetching matches for PL in 2024/2025 season...
Successfully fetched 380 matches.

Creating a single DataFrame from the collected data...
```

Woot! 6,602 unique matches is more than double what we need, and provides plenty of room to create a robust model.

What's Next: Data Cleaning and Feature Engineering

The next step is to clean and prepare this data. Our DataFrame still contains nested JSON objects in columns like `area`, `competition`, `homeTeam`, `awayTeam`, `score`, and `referees`. So we'll need "flatten" these columns to turn them into simple, usable features.

JSON is like having a Microsoft Excel sheet and you use two or three rows for miscellaneous data for each particular company/product/thing. However, it'd be a lot faster and easier by just "flattening it", where you're putting all that excess data in multiple columns on the same single row for each thing/company/etc.

Flatten the `homeTeam`, `awayTeam`, and `score` columns

We can use the `pandas.json_normalize` function to flatten the nested data into new

columns. The key information we want to focus on is:

- **Team Names:** The `name` from both the `homeTeam` and `awayTeam` objects.
 - **Final Scores:** The `fullTime` and `halfTime` scores from the `score` object.
 - **Winner:** The `winner` field from the `score` object. Obviously, this is the main target and a crucial feature for this model.
-

The Main Gooooooooooooooal Here (yeah, we had to.)

- **"How often the home side wins?"** Through a process of **descriptive analysis** on the data, we'll simply calculate the percentage of home wins in our 6,000+ matches.
- **"Does a yellow card means a loss?"** We will use features like the number of yellow/red cards to learn the **correlation** between card discipline and match outcomes. For example, like turnover margin in the NFL usually decides the winner, soccer matches may follow a similar pattern where losing teams usually receive more cards.
- **"What score at halftime means a victory?"** Halftime score can be a very strong **predictive feature**, although there's always a chance of a comeback. The model will use historical data to learn which halftime scores most often result in a victory.

THIS IS AS FAR AS I GOT TO.. the rest is Gemini.. so we'll need to edit it ourselves *****

That's a fantastic outcome! You have successfully completed the first major requirement of your project: collecting a large, original dataset via an API. The `6602 unique matches` is an excellent foundation.

Now, let's move on to the next phase of your project, as outlined in your instructions: **"Complete cleaning, preprocessing, and EDA."**

Here are the actionable steps you should take with your `final_df` to fulfill the rest of that requirement:

1. **Flatten the Data (Pre-processing):** Your DataFrame contains nested JSON objects in columns like `homeTeam`, `awayTeam`, `score`, and `referees`. You need to turn this complex data into a "flat" table with simple columns.
 - **Action:** Use `pandas.json_normalize()` to flatten the data from these `nested fields`. This will give you new columns like `homeTeam.name`.

nested fields. This will give you new columns like `homeTeam.name`, `awayTeam.name`, `score.fullTime.home`, and `score.fullTime.away`.

- **Goal:** Create a simple, flat DataFrame where each row is a single match and each column is a single feature.

2. Handle Missing Values and Justify (Cleaning):

After flattening, some values might be missing or corrupted.

- **Action:** Use `df.isnull().sum()` to identify how many missing values are in each column. Then, decide how to handle them (e.g., dropping rows, or filling with a placeholder like 'Unknown').
- **Goal:** Ensure every column you plan to use for your model has no missing values and is clean.

3. Perform Exploratory Data Analysis (EDA):

This is where you get to show off your data-driven insights.

- **Action:** Create visualizations (like bar charts or histograms) to show things like:
 - The distribution of wins, losses, and draws across all leagues.
 - The average number of goals scored per match.
 - The number of matches played per season.
- **Goal:** Write a few sentences summarizing your key findings in a markdown cell.

4. Detect and Treat Outliers:

Outliers can skew your model's results, so you need to identify them.

- **Action:** For numerical features like goals, use the Interquartile Range (IQR) method to detect extreme values (e.g., a match with an unusually high number of goals).
- **Goal:** Justify in a markdown cell whether you will remove these outliers or keep them for your analysis.

5. Encode and Scale Features:

This is the final step before you can build your model.

- **Action:**
 - Use **One-Hot Encoding** on categorical features (like `homeTeam.name`, `awayTeam.name`, and the `winner` field) to turn them into numerical data.
 - Use **Standard Scaling** on any numerical features (like `goals_scored` or `yellow_cards`) to ensure all features are on the same scale.

- **Goal:** Prepare a final, clean, and processed DataFrame that's ready to be fed into your machine learning model.

```

#
#
# THIS IS THE 'BEFORE' PREPROCESSING..
#
#

import pandas as pd

# Assuming 'final_df' is the DataFrame you created in the previous step
# If not, you'll need to run the data collection script again to create it

print("Starting data cleaning and preprocessing...\n")

# --- Step 1: Flattening the Data ---
# We'll extract key information from the nested columns.

# Extract team names and IDs from homeTeam and awayTeam columns
df_teams = pd.json_normalize(final_df['homeTeam']).add_prefix('homeTeam_')
df_teams = df_teams.join(pd.json_normalize(final_df['awayTeam']).add_prefix('aw'))

# Extract score details from the 'score' column
df_scores = pd.json_normalize(final_df['score'])

# Combine all the flattened data into a new DataFrame
cleaned_df = pd.concat([
    final_df[['id', 'utcDate', 'status', 'competition']], # Keep essential cols
    df_teams,
    df_scores
], axis=1)

# --- Step 2: Feature Engineering and Dropping Unnecessary Columns ---
# Create a winner column based on the 'winner' field in the score data
cleaned_df['winner'] = cleaned_df['winner'].fillna('DRAW') # Handle matches that

# Drop the original nested columns and other irrelevant columns
columns_to_drop = ['homeTeam', 'awayTeam', 'score', 'area', 'odds', 'referees']
cleaned_df = cleaned_df.drop(columns=columns_to_drop, errors='ignore')

# --- Step 3: Check for Missing Values ---
print("Checking for any remaining missing values...\n")
missing_values = cleaned_df.isnull().sum()
print("Missing values per column:\n", missing_values[missing_values > 0])
print(f"\nDataFrame now has {len(cleaned_df.columns)} columns and {len(cleaned_df)} rows")

print("Cleaned DataFrame is ready. Here are the first 5 rows and columns:")
print(cleaned_df.head())

```

```
print("\nNew DataFrame Columns:")
print(cleaned_df.columns)
```

Starting data cleaning and preprocessing...

Checking for any remaining missing values...

Missing values per column:

```
halfTime.home      1
halfTime.away      1
regularTime.home   6601
regularTime.away   6601
extraTime.home     6601
extraTime.away     6601
dtype: int64
```

DataFrame now has 24 columns and 6602 rows.

Cleaned DataFrame is ready. Here are the first 5 rows and columns:

```
id          utcDate    status  \
0  441789  2023-08-18T18:30:00Z FINISHED
1  441792  2023-08-19T13:30:00Z FINISHED
2  441794  2023-08-19T13:30:00Z FINISHED
3  441795  2023-08-19T13:30:00Z FINISHED
4  441796  2023-08-19T13:30:00Z FINISHED

competition  homeTeam_id  \
0  {'id': 2002, 'name': 'Bundesliga', 'code': 'BL...'      12
1  {'id': 2002, 'name': 'Bundesliga', 'code': 'BL...'       3
2  {'id': 2002, 'name': 'Bundesliga', 'code': 'BL...'      11
3  {'id': 2002, 'name': 'Bundesliga', 'code': 'BL...'       2
4  {'id': 2002, 'name': 'Bundesliga', 'code': 'BL...'      16

homeTeam_name  homeTeam_shortName  homeTeam_tla  \
0  SV Werder Bremen           Bremen        SVW
1  Bayer 04 Leverkusen       Leverkusen      B04
2  VfL Wolfsburg            Wolfsburg      WOB
3  TSG 1899 Hoffenheim       Hoffenheim    TSG
4  FC Augsburg              Augsburg       FCA

homeTeam_crest  awayTeam_id  ...  winner  \
0  https://crests.football-data.org/12.png      5  ...  AWAY_TEAM
1  https://crests.football-data.org/3.png    721  ...  HOME_TEAM
2  https://crests.football-data.org/11.png     44  ...  HOME_TEAM
3  https://crests.football-data.org/2.png     17  ...  AWAY_TEAM
4  https://crests.football-data.org/16.png    18  ...  DRAW

duration fullTime.home fullTime.away halfTime.home halfTime.away  \
0  REGULAR          0            4        0.0        1.0
1  REGULAR          3            2        2.0        1.0
2  REGULAR          2            0        2.0        0.0
3  REGULAR          1            2        0.0        2.0
4  REGULAR          4            4        3.0        3.0

regularTime.home  regularTime.away  extraTime.home  extraTime.away
```

0	NaN	NaN	NaN	NaN
1	NaN	NaN	NaN	NaN
2	NaN	NaN	NaN	NaN
3	NaN	NaN	NaN	NaN
4	NaN	NaN	NaN	NaN

The output you've provided shows a great start, but it also reveals a few new challenges that we need to address. This is a very common part of the data cleaning process!

Issues to Fix:

Lots of missing values like regularTime.home, regularTime.away, extraTime.home, and extraTime.away columns. Looks like these columns only contain data for cup or tournament matches that go into extra time or penalties, so that's not really relevant for league matches.

Remaining nested data where the output of cleaned_df.head() shows that the competition column is still nested, so we'll need to flatten that to get the league name.

```

#
#
# THIS IS A BETTER 'AFTER' PREPROCESSING..
#
#

import pandas as pd

# Assuming 'final_df' is the DataFrame you created earlier

print("Starting a more comprehensive data cleaning and preprocessing...\n")

# --- Step 1: Flattening the Data ---
# Extract key information from the nested columns.

# Normalize home and away team data
home_teams = pd.json_normalize(final_df['homeTeam']).add_prefix('homeTeam_')
away_teams = pd.json_normalize(final_df['awayTeam']).add_prefix('awayTeam_')

# Normalize score data
scores = pd.json_normalize(final_df['score'])

# Normalize competition data to get the league name
competitions = pd.json_normalize(final_df['competition']).add_prefix('competiti')

# Combine all the flattened data into a new DataFrame
cleaned_df = pd.concat([
    final_df[['id', 'utcDate', 'status']], # Keep essential columns
    competitions[['competition_name', 'competition_code']],
    ...
])

```

```

    home_teams[['homeTeam_name', 'homeTeam_tla']],
    away_teams[['awayTeam_name', 'awayTeam_tla']],
    scores[['winner', 'fullTime.home', 'fullTime.away', 'halfTime.home', 'halfTime.away'],
], axis=1)

# --- Step 2: Handle Missing Values ---
# Fill the few missing halftime scores with 0
cleaned_df['halfTime.home'] = cleaned_df['halfTime.home'].fillna(0).astype(int)
cleaned_df['halfTime.away'] = cleaned_df['halfTime.away'].fillna(0).astype(int)

# Drop columns with too many missing values
cleaned_df.drop(columns=['regularTime.home', 'regularTime.away', 'extraTime.home',
                        'extraTime.away'], inplace=True)

# --- Step 3: Rename Columns for Clarity ---
cleaned_df.rename(columns={
    'fullTime.home': 'home_score',
    'fullTime.away': 'away_score',
    'halfTime.home': 'halfTime_home_score',
    'halfTime.away': 'halfTime_away_score'
}, inplace=True)

print("Comprehensive cleaning complete.\n")
print(f"DataFrame now has {len(cleaned_df.columns)} columns and {len(cleaned_df)} rows")
print("The first 5 rows of the new, clean DataFrame:")
print(cleaned_df.head())
print("\nNew DataFrame Columns:")
print(cleaned_df.columns)

```

Starting a more comprehensive data cleaning and preprocessing...

Comprehensive cleaning complete.

DataFrame now has 14 columns and 6602 rows.

The first 5 rows of the new, clean DataFrame:

	id	utcDate	status	competition_name	competition_code	
0	441789	2023-08-18T18:30:00Z	FINISHED	Bundesliga	BL1	
1	441792	2023-08-19T13:30:00Z	FINISHED	Bundesliga	BL1	
2	441794	2023-08-19T13:30:00Z	FINISHED	Bundesliga	BL1	
3	441795	2023-08-19T13:30:00Z	FINISHED	Bundesliga	BL1	
4	441796	2023-08-19T13:30:00Z	FINISHED	Bundesliga	BL1	
	homeTeam_name	homeTeam_tla		awayTeam_name	awayTeam_tla	
0	SV Werder Bremen	SVW		FC Bayern München	FCB	
1	Bayer 04 Leverkusen	B04		RB Leipzig	RBL	
2	VfL Wolfsburg	WOB	1. FC Heidenheim 1846		HEI	
3	TSG 1899 Hoffenheim	TSG		SC Freiburg	SCF	
4	FC Augsburg	FCA	Borussia Mönchengladbach		BMG	
	winner	home_score	away_score	halfTime_home_score	halfTime_away_score	
0	AWAY_TEAM	0	4	0		1
1	HOME_TEAM	3	2	2		1
2	HOME_TEAM	2	0	2		0
3	AWAY_TEAM	1	2	0		2

```
4      DRAW      4      4      3      3
```

```
New DataFrame Columns:  
Index(['id', 'utcDate', 'status', 'competition_name', 'competition_code',  
       'homeTeam_name', 'homeTeam_tla', 'awayTeam_name', 'awayTeam_tla',  
       'winner', 'home_score', 'away_score', 'halfTime_home_score',  
       'halfTime_away_score'],  
      dtype='object')
```

```
import matplotlib.pyplot as plt  
import seaborn as sns  
  
# Make sure you have the 'cleaned_df' from the previous step loaded in your notebook  
  
# Set a clean visual style for the plots  
sns.set_style("whitegrid")  
  
print("Performing Exploratory Data Analysis...\n")  
  
# --- 1. Overall Match Outcome Distribution ---  
# This answers your question: "how often the home side wins?"  
plt.figure(figsize=(8, 6))  
sns.countplot(x='winner', data=cleaned_df, palette='viridis')  
plt.title('Overall Match Outcome Distribution (All Leagues)', fontsize=16)  
plt.xlabel('Winner', fontsize=12)  
plt.ylabel('Number of Matches', fontsize=12)  
plt.show()  
  
# --- 2. Score Distribution ---  
# This shows the distribution of goals scored by home and away teams.  
plt.figure(figsize=(12, 6))  
sns.histplot(cleaned_df['home_score'], bins=range(0, 10), kde=False, color='skyblue')  
sns.histplot(cleaned_df['away_score'], bins=range(0, 10), kde=False, color='salmon')  
plt.title('Distribution of Goals Scored (Full Time)', fontsize=16)  
plt.xlabel('Number of Goals', fontsize=12)  
plt.ylabel('Number of Matches', fontsize=12)  
plt.xticks(range(10))  
plt.legend()  
plt.show()  
  
# --- 3. Match Outcomes by League ---  
# This helps you compare the win rates between different leagues.  
plt.figure(figsize=(15, 8))  
sns.countplot(x='competition_code', hue='winner', data=cleaned_df, palette='viridis')  
plt.title('Match Outcomes by League', fontsize=16)  
plt.xlabel('League Code', fontsize=12)  
plt.ylabel('Number of Matches', fontsize=12)  
plt.legend(title='Winner')
```

```

plt.xticks(rotation=45)
plt.show()

print("Calculating and plotting match outcome percentages by league...\n")

# --- Step 1: Calculate the total number of matches per league ---
total_matches_per_league = cleaned_df.groupby('competition_code')['id'].count()

# --- Step 2: Calculate the count for each outcome per league ---
outcome_counts = cleaned_df.groupby(['competition_code', 'winner']).size().unstack()

# --- Step 3: Calculate the percentages ---
outcome_percentages = outcome_counts.div(total_matches_per_league, axis=0) * 100

# --- Step 4: Plotting the results ---
outcome_percentages.plot(kind='bar', figsize=(15, 8), color=['skyblue', 'salmon'])

plt.title('Match Outcome Percentages by League', fontsize=16)
plt.xlabel('League Code', fontsize=12)
plt.ylabel('Percentage of Matches (%)', fontsize=12)
plt.xticks(rotation=45)
plt.legend(title='Outcome', loc='upper right')
plt.tight_layout()
plt.show()

print("Calculating and formatting the new table with Home/Away win difference..")

# --- Step 1: Recalculate the percentages ---
# Calculate total matches per league
total_matches_per_league = cleaned_df.groupby('competition_code')['id'].count()

# Count the number of matches for each outcome per league
outcome_counts = cleaned_df.groupby(['competition_code', 'winner']).size().unstack()

# Calculate the percentages for each outcome and round to 2 decimals
outcome_percentages = (outcome_counts.div(total_matches_per_league, axis=0) * 100).round(2)

# --- Step 2: Add the new column for Home/Away win difference ---
outcome_percentages['Home_Away_Win_Difference'] = outcome_percentages['HOME_TEAM_WIN_PCT'] - outcome_percentages['AWAY_TEAM_WIN_PCT']

# --- Step 3: Sort the table for better readability and print ---
# We'll sort by the new difference column to easily see the leagues with the most extreme differences
formatted_table = outcome_percentages.sort_values(by='Home_Away_Win_Difference', ascending=False)

print("Here is the formatted table, sorted by the difference in home and away wins")
print(formatted_table.to_string())

```


well.. the home team certainly seems to have an advantage, of course, like in most sports.

of course the goal differential is quite telling.. the more goals you score, the more you win!

.....

but this is quite interesting.. the parity of some leagues - Ligue 1 (FRA) and Die Bundesliga (GER) seem to have teams matched pretty well the past two seasons, where the home/away ratio is quite close. The second league in England, however, has more teams and also a wide berth of good and bad teams, so that ratio seems a lot bigger.

But whoa.. in Brazil, a 20% difference is huge.
Especially if you're betting on games. You'll probably win a lot if you bet on "Home Team or Tie" on most matches.

```
import pandas as pd

# This code assumes 'cleaned_df' is the DataFrame from the previous steps.

print("Starting outlier detection using the IQR method...\n")

# Identify the columns with numerical scores
score_columns = ['home_score', 'away_score']

# Create a new DataFrame to store outliers
outliers = pd.DataFrame()

for col in score_columns:
    Q1 = cleaned_df[col].quantile(0.25)
    Q3 = cleaned_df[col].quantile(0.75)
    IQR = Q3 - Q1

    # Define outlier bounds
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR

    # Find outliers for the current column
    col_outliers = cleaned_df[(cleaned_df[col] < lower_bound) | (cleaned_df[col] > upper_bound)]

    # Add the found outliers to our outlier DataFrame
    outliers = pd.concat([outliers, col_outliers])

    print(f"For '{col}':")
    print(f"  Q1: {Q1}, Q3: {Q3}, IQR: {IQR}")
    print(f"  Outlier Bounds: ({lower_bound}, {upper_bound})")
    print(f"  Number of outliers found: {len(col_outliers)}")

# Remove duplicates from the combined outlier DataFrame
```

```

outliers.drop_duplicates(inplace=True)

print(f"\nTotal unique matches identified as outliers: {len(outliers)}")
print("\nHere are a few examples of the matches with outlier scores:")
print(outliers.head())

```

Starting outlier detection using the IQR method...

```

For 'home_score':
    Q1: 1.0, Q3: 2.0, IQR: 1.0
    Outlier Bounds: (-0.5, 3.5)
    Number of outliers found: 507
For 'away_score':
    Q1: 0.0, Q3: 2.0, IQR: 2.0
    Outlier Bounds: (-3.0, 5.0)
    Number of outliers found: 25

```

Total unique matches identified as outliers: 531

Here are a few examples of the matches with outlier scores:

	<code>id</code>	<code>utcDate</code>	<code>status</code>	<code>competition_name</code>	<code>competition_code</code>	
4	441796	2023-08-19T13:30:00Z	FINISHED	Bundesliga	BL1	
5	441797	2023-08-19T13:30:00Z	FINISHED	Bundesliga	BL1	
7	441791	2023-08-20T13:30:00Z	FINISHED	Bundesliga	BL1	
9	441799	2023-08-25T18:30:00Z	FINISHED	Bundesliga	BL1	
19	441809	2023-09-02T13:30:00Z	FINISHED	Bundesliga	BL1	
		<code>homeTeam_name</code>	<code>homeTeam_tla</code>	<code>awayTeam_name</code>	<code>awayTeam_tla</code>	\
4		FC Augsburg	FCA	Borussia Mönchengladbach	BMG	
5		VfB Stuttgart	VFB	VfL Bochum 1848	BOC	
7		1. FC Union Berlin	UNB	1. FSV Mainz 05	M05	
9		RB Leipzig	RBL	VfB Stuttgart	VFB	
19		Bayer 04 Leverkusen	B04	SV Darmstadt 98	SVD	
		<code>winner</code>	<code>home_score</code>	<code>away_score</code>	<code>halfTime_home_score</code>	\
4		DRAW	4	4	3	
5		HOME_TEAM	5	0	2	
7		HOME_TEAM	4	1	2	
9		HOME_TEAM	5	1	0	
19		HOME_TEAM	5	1	1	
			<code>halfTime_away_score</code>			
4			3			
5			0			
7			0			
9			1			
19			1			

- ✓ no.. these outliers stay!

these are part of the game. high scores happen.

```

import pandas as pd
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.compose import ColumnTransformer

# This code assumes 'cleaned_df' is the DataFrame from the previous steps.

print("Starting final feature encoding and scaling...\n")

# --- Separate features (X) from the target variable (y) ---
# We explicitly select ONLY the features we need
categorical_features = ['competition_code', 'homeTeam_tla', 'awayTeam_tla']
numerical_features = ['home_score', 'away_score', 'halfTime_home_score', 'halfT

# Create a new DataFrame with ONLY the columns you will use as features
X = cleaned_df[categorical_features + numerical_features].copy()

# The target variable (y) is the 'winner' column
y = cleaned_df['winner']

# --- Set up the preprocessor ---
preprocessor = ColumnTransformer(
    transformers=[
        ('num', StandardScaler(), numerical_features),
        ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features)
    ],
    # 'remainder' is not needed since we've already selected our features
)

# --- Apply the transformations to the features ---
X_transformed = preprocessor.fit_transform(X)

# --- One-hot encode the target variable ('winner') ---
ohe_winner = OneHotEncoder(handle_unknown='ignore')
y_encoded = ohe_winner.fit_transform(y.values.reshape(-1, 1))

print("Features have been successfully encoded and scaled.")
print(f"Final feature matrix (X) shape: {X_transformed.shape}")
print(f"Final target variable (y) shape: {y_encoded.shape}")
print("\nYour data is now ready for model training! 🎉")

```

Starting final feature encoding and scaling...

Features have been successfully encoded and scaled.
 Final feature matrix (X) shape: (6602, 391)
 Final target variable (y) shape: (6602, 3)

Your data is now ready for model training! 🎉

Engineering New Features for the Soccer Data

Now that our soccer data is clean, let's create some new pieces of information from the existing ones. Think of this like taking ingredients you already have (like flour and sugar) and combining them to make something new (like a cake!). These new pieces of information, or "features," can help our model understand the game better.

Feature 1: Goal Difference

This is a simple but important one. We'll figure out the difference between the number of goals the home team scored and the number of goals the away team scored.

Why is this helpful? The goal difference tells us how much one team dominated the other in terms of scoring. A big positive number means the home team won comfortably, while a big negative number means the away team won easily. This is a direct measure of how one-sided a match was.

```
# Calculate the goal difference for each match
cleaned_df['goal_difference'] = cleaned_df['home_score'] - cleaned_df['away_score']

print("Added 'goal_difference' column.")
display(cleaned_df[['home_score', 'away_score', 'goal_difference']].head())
```

▼ Feature 2: Total Goals

This feature is just the total number of goals scored in a match by both teams combined.

Why is this helpful? The total number of goals can give us an idea of the style of play in a match. A high total might indicate an open, attacking game, while a low total could suggest a more defensive or cautious approach. This could be useful for understanding match dynamics beyond just who won.

```
# Calculate the total number of goals in each match
cleaned_df['total_goals'] = cleaned_df['home_score'] + cleaned_df['away_score']
```

```
print("\nAdded 'total_goals' column.")
display(cleaned_df[['home_score', 'away_score', 'total_goals']].head())
```

▼ Feature 3: Extracting Date and Time Information

The original data has a timestamp for when the match happened. We can pull out different parts of this timestamp, like the hour of the day, the day of the week, or the month.

Why is this helpful? This could capture patterns related to scheduling or time of year. For example, games played on weekends might have different dynamics than weekday games, or teams might perform differently at certain times of the year due to fatigue or weather.

```
# Convert the 'utcDate' column to datetime objects
cleaned_df['utcDate'] = pd.to_datetime(cleaned_df['utcDate'])

# Extract date and time components
cleaned_df['match_hour'] = cleaned_df['utcDate'].dt.hour
cleaned_df['match_day_of_week'] = cleaned_df['utcDate'].dt.dayofweek # Monday=0
cleaned_df['match_month'] = cleaned_df['utcDate'].dt.month

print("\nAdded date and time features.")
display(cleaned_df[['utcDate', 'match_hour', 'match_day_of_week', 'match_month']])
```

▼ Feature 4: Home Advantage Indicator

This is a simple "yes" or "no" feature that just tells us if the match was played at the home team's stadium.

Why is this helpful? As we saw in our earlier analysis, playing at home often provides an advantage due to familiar surroundings, crowd support, and less travel fatigue. This feature directly captures that important factor.

```
# Create a binary feature for home advantage (1 if home game, 0 otherwise)
# This is technically captured by the structure of the data and how we'll split
# but having an explicit column can sometimes be useful depending on the model.
# For a simple classification model predicting winner, this isn't strictly nece
# will learn the home/away relationship from the homeTeam/awayTeam features,
# but it's included as an example of a simple binary feature.
cleaned_df['is_home_game'] = 1 # All rows in this DataFrame represent a game fr
                                # Since we are analyzing from the perspective of t
                                # this feature would always be 1.
                                # A more useful "home advantage" feature might be
                                # to have each row be a team's performance in a ma

# Given the current DataFrame structure where each row is a match with home and
# the 'winner' column already implicitly captures home advantage.
# Let's create a feature that might be more directly interpretable or useful:
# A binary feature indicating if the *home team* won.

cleaned_df['home_team_won'] = (cleaned_df['winner'] == 'HOME_TEAM').astype(int)
cleaned_df['away_team_won'] = (cleaned_df['winner'] == 'AWAY_TEAM').astype(int)
cleaned_df['is_draw'] = (cleaned_df['winner'] == 'DRAW').astype(int)

print("\nAdded home advantage and outcome indicator features.")
display(cleaned_df[['winner', 'home_team_won', 'away_team_won', 'is_draw']]).hea
```

```

# --- Create safe pre-match proxy for goal_difference --- TESTING CODE!!!!!!!
# --- Create safe pre-match proxy for goal_difference ---

# Calculate per-match goal differences (post-match, but we'll use them only for
cleaned_df['home_goal_diff'] = cleaned_df['home_score'] - cleaned_df['away_score']
cleaned_df['away_goal_diff'] = cleaned_df['away_score'] - cleaned_df['home_score']

# Sort by date so rolling windows are correct
cleaned_df = cleaned_df.sort_values(by='utcDate')

# Rolling averages for last 5 matches (excluding current match)
cleaned_df['home_avg_goal_diff_last5'] = (
    cleaned_df.groupby('homeTeam_name')['home_goal_diff']
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)

cleaned_df['away_avg_goal_diff_last5'] = (
    cleaned_df.groupby('awayTeam_name')['away_goal_diff']
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)

# Combine into a single "expected goal diff" feature
cleaned_df['expected_goal_diff'] = (
    cleaned_df['home_avg_goal_diff_last5'] - cleaned_df['away_avg_goal_diff_last5']
)

```

```

# TESTING TESTING 7+$##%#(*%#^^^^^^^^^^^^^^^^^^^^^^^^^)
#
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
from sklearn.preprocessing import OneHotEncoder
import pandas as pd

# Remove the known leaky feature(s)
X_no_leak = X.drop(columns=['goal_difference'], errors='ignore')

print("Testing each feature individually (goal_difference removed):\n")

for col in X_no_leak.columns:
    X_col = X_no_leak[[col]].copy()

    # Drop rows with NaN in this column
    mask = X_col[col].notna()
    X_col = X_col.loc[mask]

```

```
x_col = x_col[mask]
y_col = y[mask]

# One-hot encode categorical features
if X_col[col].dtype == 'object' or str(X_col[col].dtype).startswith('categc
    enc = OneHotEncoder(handle_unknown='ignore', sparse_output=False)
    X_col = pd.DataFrame(enc.fit_transform(X_col))

# Train/test split
X_train, X_test, y_train, y_test = train_test_split(
    X_col, y_col, test_size=0.2, random_state=42, stratify=y_col
)

# Train and evaluate
model = LogisticRegression(max_iter=1000)
model.fit(X_train, y_train)
acc = accuracy_score(y_test, model.predict(X_test))
print(f'{col}: {acc:.4f}')
```

Testing each feature individually (goal_difference removed):

```
competition_code: 0.4375  
homeTeam_tla: 0.4504  
awayTeam_tla: 0.4565  
home_score: 0.5942  
away_score: 0.6268  
halfTime_home_score: 0.5201  
halfTime_away_score: 0.5360
```

DATA LEAKAGE ISSUE!

As indicated by the code comments above, we noticed the accuracy and recall and everything was *perfect*. 100%. Wonderful!

Which is basically wishful thinking.. there had to be an issue somewhere. After running some tests, **we realized the 'goal differential' feature was a problem - because it was predicting the winner AFTER the match was over.**

Obviously it'll guess the winner!

So we stopped using that feature.. however, we made a feature to determine the goal diff for the *previous* 5 matches, not the match being played that moment.

```

from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.compose import ColumnTransformer
import pandas as pd

# -----
# 1. Remove known leaky feature(s)
# -----
X_no_leak = X.drop(columns=['goal_difference'], errors='ignore')

print("== Leakage Check: Single Feature Accuracy ==")
for col in X_no_leak.columns:
    X_col = X_no_leak[[col]].copy()

    # Drop rows with NaN in this column
    mask = X_col[col].notna()
    X_col = X_col[mask]
    y_col = y[mask]

    # One-hot encode categorical features
    if X_col[col].dtype == 'object' or str(X_col[col].dtype).startswith('categc'):
        enc = OneHotEncoder(handle_unknown='ignore', sparse_output=False)
        X_col = pd.DataFrame(enc.fit_transform(X_col))

    # Train/test split
    X_train, X_test, y_train, y_test = train_test_split(
        X_col, y_col, test_size=0.2, random_state=42, stratify=y_col
    )

    # Train and evaluate
    model = LogisticRegression(max_iter=1000)
    model.fit(X_train, y_train)
    acc = accuracy_score(y_test, model.predict(X_test))
    print(f"{col}: {acc:.4f}")

# -----
# 2. Define cleaned feature lists
# -----
categorical_features = [
    'competition_name', 'competition_code',
    'homeTeam_name', 'homeTeam_tla',
    'awayTeam_name', 'awayTeam_tla'
]

# Replace 'goal_difference' with safe proxy if available
numerical_features = [
    'expected_goal_diff', # <-- safe pre-match proxy (must exist in cleaned_df
    'match_hour', 'match_day_of_week', 'match_month',
    'is_home_game',
    'homeTeam_rolling_avg_goals_scored_5'
]

```

```

        'homeTeam_rolling_avg_goals_scored_5' ,
        'awayTeam_rolling_avg_goals_conceded_5'
    ]

# Filter out any features that don't exist in the DataFrame
numerical_features = [f for f in numerical_features if f in cleaned_df.columns]

# -----
# 3. Split-first preprocessing
# -----
X_clean = cleaned_df[categorical_features + numerical_features]
y_clean = cleaned_df['winner']

X_train, X_test, y_train, y_test = train_test_split(
    X_clean, y_clean, test_size=0.2, random_state=42, stratify=y_clean
)

from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.compose import ColumnTransformer

# Preprocessor with imputation for numeric features
preprocessor = ColumnTransformer(
    transformers=[
        ('num', Pipeline([
            ('imputer', SimpleImputer(strategy='constant', fill_value=0)), # f
            ('scaler', StandardScaler())
        ]), numerical_features),

        ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features)
    ]
)

X_train_processed = preprocessor.fit_transform(X_train)
X_test_processed = preprocessor.transform(X_test)

# -----
# 4. Train Logistic Regression
# -----
log_reg = LogisticRegression(max_iter=1000)
log_reg.fit(X_train_processed, y_train)

y_pred = log_reg.predict(X_test_processed)

print("\n==== Final Logistic Regression Results (Leakage-Free) ===")
print("Accuracy:", accuracy_score(y_test, y_pred))
print("\nClassification Report:\n", classification_report(y_test, y_pred))

```

```

==== Leakage Check: Single Feature Accuracy ====
competition_code: 0.4375
homeTeam_tla: 0.4504
awayTeam_tla: 0.4565
home_score: 0.5942
away_score: 0.6268
halfTime_home_score: 0.5201
halfTime_away_score: 0.5360

==== Final Logistic Regression Results (Leakage-Free) ====
Accuracy: 0.5185465556396669

```

Classification Report:

	precision	recall	f1-score	support
AWAY_TEAM	0.50	0.50	0.50	404
DRAW	0.30	0.12	0.17	339
HOME_TEAM	0.56	0.76	0.65	578
accuracy			0.52	1321
macro avg	0.45	0.46	0.44	1321
weighted avg	0.48	0.52	0.48	1321

```

from sklearn.model_selection import train_test_split

# Assuming X_processed contains your processed features and y contains your target
# X_processed was created in the previous steps by encoding and scaling the features
# y is the 'winner' column from cleaned_df

# Split the data into training and testing sets
# We'll use a common split ratio, like 80% for training and 20% for testing
# Setting a random_state ensures that the split is the same each time you run it
X_train, X_test, y_train, y_test = train_test_split(
    X_transformed, y, test_size=0.2, random_state=42, stratify=y # Stratify to keep proportions
)

print("Data successfully split into training and testing sets.")
print(f"Training features shape: {X_train.shape}")
print(f"Testing features shape: {X_test.shape}")
print(f"Training target shape: {y_train.shape}")
print(f"Testing target shape: {y_test.shape}")

Data successfully split into training and testing sets.
Training features shape: (5281, 391)
Testing features shape: (1321, 391)
Training target shape: (5281,)
Testing target shape: (1321,)

```

✓ Training a Classical Machine Learning Model: Logistic Regression

Now that our soccer data is ready, let's use a classical machine learning method called **Logistic Regression** to build our first model. Think of this model as trying to find the best way to draw lines (or more complex boundaries) in our data to separate the different outcomes (Home Win, Away Win, Draw).

```
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report

# We already have our data split into training sets (X_train, y_train)
# and testing sets (X_test, y_test) from the previous step.

print("Let's train our Logistic Regression model...")

# Here, we're setting up the Logistic Regression model.
# max_iter is like telling the model to take enough steps to learn properly.
# random_state makes sure we get the same results each time we run this code.
model = LogisticRegression(max_iter=1000, random_state=42)

# This is where the model learns from the training data.
# It looks at the features (X_train) and the actual outcomes (y_train)
# to find patterns.
model.fit(X_train, y_train)

print("Model training complete! The model has learned from the data.")

# Now, let's see how well our trained model performs on data it hasn't seen before.
# We ask the model to predict the outcomes (y_pred) based on the features in the test data.
y_pred = model.predict(X_test)

print("\nNow let's evaluate how good our model is at predicting outcomes on the test data.")

# Accuracy tells us the overall percentage of predictions the model got right.
accuracy = accuracy_score(y_test, y_pred)
print(f"Overall Accuracy: {accuracy:.4f}\n")

# The Classification Report gives us more detailed information about the model's performance
# for each possible outcome (Home Win, Away Win, Draw).
print("Detailed Performance Report (Classification Report):")
print(classification_report(y_test, y_pred))

Let's train our Logistic Regression model...
Model training complete! The model has learned from the data.

Now let's evaluate how good our model is at predicting outcomes on the test data
Overall Accuracy: 1.0000

Detailed Performance Report (Classification Report):
precision    recall    f1-score    support
AWAY  TFAM      1.00      1.00      1.00      404

```

	1.00	1.00	1.00	TOT
DRAW	1.00	1.00	1.00	339
HOME_TEAM	1.00	1.00	1.00	578
accuracy			1.00	1321
macro avg	1.00	1.00	1.00	1321
weighted avg	1.00	1.00	1.00	1321

```

from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report

# Assuming X_processed and y are defined from the previous steps

# Split the data into training and testing sets again with the corrected X_proc
# We need to re-split because X_processed has changed
X_train, X_test, y_train, y_test = train_test_split(
    X_transformed, y, test_size=0.2, random_state=42, stratify=y
)

print("Data successfully re-split into training and testing sets with corrected")
print(f"Training features shape: {X_train.shape}")
print(f"Testing features shape: {X_test.shape}")
print(f"Training target shape: {y_train.shape}")
print(f"Testing target shape: {y_test.shape}")

print("\nRetraining Logistic Regression model with corrected features...\n")

# Initialize the Logistic Regression model
model = LogisticRegression(max_iter=1000, random_state=42)

# Train the model using the training data
model.fit(X_train, y_train)

print("Model retraining complete.")

# Make predictions on the testing data
y_pred = model.predict(X_test)

print("\nModel evaluation on the test set with corrected features:")

# Evaluate the model's accuracy
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy:.4f}\n")

# Print a classification report for more detailed metrics
print("Classification Report:")
print(classification_report(y_test, y_pred))

```

```
Data successfully re-split into training and testing sets with corrected feature  
Training features shape: (5281, 391)  
Testing features shape: (1321, 391)  
Training target shape: (5281,)  
Testing target shape: (1321,)
```

```
Retraining Logistic Regression model with corrected features...
```

```
Model retraining complete.
```

```
Model evaluation on the test set with corrected features:  
Accuracy: 1.0000
```

```
Classification Report:
```

	precision	recall	f1-score	support
AWAY_TEAM	1.00	1.00	1.00	404
DRAW	1.00	1.00	1.00	339
HOME_TEAM	1.00	1.00	1.00	578
accuracy			1.00	1321
macro avg	1.00	1.00	1.00	1321
weighted avg	1.00	1.00	1.00	1321

▼ Training Another Classical Model: Random Forest

Alright team, let's move on to trying another way to predict our soccer match outcomes. Remember how we used Logistic Regression? That was like using one straightforward method to make predictions.

Now, we're going to try something called a **Random Forest Classifier**. Think of it like this: instead of just one of us making a prediction, we're going to ask a whole bunch of our friends (these are the "trees" in the "forest") to *each* make a prediction based on slightly different parts of the data. Then, we'll all get together and figure out the best overall prediction by seeing what most of our friends predicted.

Why are we doing this? Well, sometimes having a "crowd" of predictors can be more accurate than just relying on one. We want to see if this Random Forest approach can do a better job predicting match outcomes than our Logistic Regression model did. This is part of our plan to compare different types of machine learning methods for our project.

```
from sklearn.ensemble import RandomForestClassifier  
from sklearn.metrics import accuracy_score, classification_report  
  
# We'll use the same split data (X_train, X_test, y_train, y_test) that we prep  
  
print("Okay team, let's train our Random Forest model now...")
```

```

# This line sets up our Random Forest.
# n_estimators=100 means we're using 100 "trees" (our prediction friends).
# random_state=42 just helps us get the same result each time we run this code.
model_rf = RandomForestClassifier(n_estimators=100, random_state=42)

# This is where we "train" our Random Forest.
# The model looks at the training data (X_train) and the actual outcomes (y_train)
# to learn how to make predictions.
model_rf.fit(X_train, y_train)

print("Random Forest model training complete! Our prediction friends have learned!")

# Now, let's see how well our trained model (our group of friends) does on data we haven't seen before.
# We ask the model to predict the outcomes (y_pred_rf) based on the features in X_test.
y_pred_rf = model_rf.predict(X_test)

print("\nLet's check how well our Random Forest model predicted the test matches.")

# Accuracy tells us the overall percentage of predictions the model got right.
accuracy_rf = accuracy_score(y_test, y_pred_rf)
print(f"Overall Accuracy: {accuracy_rf:.4f}\n")

# This Classification Report gives us more detailed scores for predicting Home vs. Away teams.
print("Detailed Performance Report (Classification Report):")
print(classification_report(y_test, y_pred_rf))

```

Okay team, let's train our Random Forest model now...
 Random Forest model training complete! Our prediction friends have learned.

Let's check how well our Random Forest model predicted the test matches.
 Overall Accuracy: 0.9939

	precision	recall	f1-score	support
AWAY_TEAM	1.00	1.00	1.00	404
DRAW	0.99	0.99	0.99	339
HOME_TEAM	0.99	1.00	0.99	578
accuracy			0.99	1321
macro avg	0.99	0.99	0.99	1321
weighted avg	0.99	0.99	0.99	1321

▼ DATA LEAKAGE FIXED

We wanted to show how one RF model that used goal diff was perfect.

Yet, when we fixed the data leakage issue, it crashed down and basically became

realistic.

Not great for the results.. but we can now work with it!

And we ran through one quick runthrough of several different regression techniques to make sure the data leakage is gone.. and see if we can tweak something to improve it.

```
#  
!pip install xgboost  
# CHECKING MULTIPLE WAYS  
  
Requirement already satisfied: xgboost in /usr/local/lib/python3.12/dist-packages  
Requirement already satisfied: numpy in /usr/local/lib/python3.12/dist-packages  
Requirement already satisfied: nvidia-nccl-cu12 in /usr/local/lib/python3.12/dist-packages  
Requirement already satisfied: scipy in /usr/local/lib/python3.12/dist-packages  
  
# Classification Not Regression *****  
#  
# LET'S PUT IT INTO ONE FINAL PROCESS - TO COMPARE AND ENSURE NO LEAKAGE  
#  
#  
import pandas as pd  
import numpy as np  
from sklearn.model_selection import train_test_split  
from sklearn.linear_model import LogisticRegression  
from sklearn.metrics import accuracy_score, classification_report  
from sklearn.preprocessing import OneHotEncoder, StandardScaler, LabelEncoder  
from sklearn.compose import ColumnTransformer  
from sklearn.svm import SVC  
from sklearn.ensemble import GradientBoostingClassifier, RandomForestClassifier  
from sklearn.neighbors import KNeighborsClassifier  
from xgboost import XGBClassifier  
from sklearn.pipeline import Pipeline  
from sklearn.impute import SimpleImputer  
  
print("Starting a comprehensive data cleaning and modeling process...\n")  
  
# --- Step 1: Data Gathering (API Call and Initial DataFrame Creation) ---  
# NOTE: This assumes you have already run the API call and have the 'final_df'  
# If you are starting from scratch, you must include your API call script here.  
  
# --- Step 2: Comprehensive Data Cleaning and Feature Engineering ---  
# Flatten nested JSON and rename columns  
home_teams = pd.json_normalize(final_df['homeTeam']).add_prefix('homeTeam_')  
away_teams = pd.json_normalize(final_df['awayTeam']).add_prefix('awayTeam_')  
scores = pd.json_normalize(final_df['score'])  
competitions = pd.json_normalize(final_df['competition']).add_prefix('competition_')  
  
cleaned_df = pd.concat([
```

```

final_df[['id', 'utcDate', 'status', 'season']],
competitions[['competition_name', 'competition_code']],
home_teams[['homeTeam_name', 'homeTeam_tla']],
away_teams[['awayTeam_name', 'awayTeam_tla']],
scores[['winner', 'fullTime.home', 'fullTime.away', 'halfTime.home', 'halfTime.away']],
], axis=1)

# Handle missing values and rename columns
cleaned_df['halfTime.home'] = cleaned_df['halfTime.home'].fillna(0).astype(int)
cleaned_df['halfTime.away'] = cleaned_df['halfTime.away'].fillna(0).astype(int)
cleaned_df.rename(columns={
    'fullTime.home': 'home_score',
    'fullTime.away': 'away_score',
    'halfTime.home': 'halfTime_home_score',
    'halfTime.away': 'halfTime_away_score'
}, inplace=True)
cleaned_df['winner'] = cleaned_df['winner'].fillna('DRAW')
columns_to_drop = ['homeTeam', 'awayTeam', 'score', 'competition', 'area', 'odd']
cleaned_df = cleaned_df.drop(columns=columns_to_drop, errors='ignore')

# Create safe pre-match proxy features
cleaned_df['home_goal_diff'] = cleaned_df['home_score'] - cleaned_df['away_score']
cleaned_df['away_goal_diff'] = cleaned_df['away_score'] - cleaned_df['home_score']
cleaned_df = cleaned_df.sort_values(by='utcDate')

cleaned_df['home_avg_goal_diff_last5'] = (
    cleaned_df.groupby('homeTeam_name')['home_goal_diff']
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)
cleaned_df['away_avg_goal_diff_last5'] = (
    cleaned_df.groupby('awayTeam_name')['away_goal_diff']
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)
cleaned_df['expected_goal_diff'] = (
    cleaned_df['home_avg_goal_diff_last5'] - cleaned_df['away_avg_goal_diff_last5']
)

print("Comprehensive cleaning and feature engineering complete.\n")

# --- Step 3: Define Features and Split Data ---
# Select only the features that are known BEFORE the match starts
categorical_features = ['competition_name', 'homeTeam_name', 'awayTeam_name']
numerical_features = ['expected_goal_diff', 'home_avg_goal_diff_last5', 'away_avg_goal_diff_last5']

# Ensure the selected features exist
numerical_features = [f for f in numerical_features if f in cleaned_df.columns]
X = cleaned_df[categorical_features + numerical_features]
y = cleaned_df['winner']

X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42, stratify=y
)

```

```

        X, y, test_size=0.2, random_state=42, stratify=y
    )

print("Data successfully split into training and testing sets.")

# --- Step 4: Preprocess Data with ColumnTransformer and Imputation ---
# A SimpleImputer is added to handle the NaN values before scaling
numerical_pipeline = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='median')),
    ('scaler', StandardScaler())
])

preprocessor = ColumnTransformer(
    transformers=[
        ('num', numerical_pipeline, numerical_features),
        ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features)
    ],
    remainder='passthrough'
)

X_processed = preprocessor.fit_transform(X) # Use X here instead of X_train
print("Features processed successfully with the specified features.\n")
print(f"Shape of X_processed: {X_processed.shape}")

# --- Step 5: Train and Compare Models ---
# Encode labels for XGBoost
le = LabelEncoder()
y_enc = le.fit_transform(y) # Encode all y

# List of CLASSIFICATION models to test
models = [
    ("Logistic Regression", LogisticRegression(max_iter=1000)),
    ("Random Forest", RandomForestClassifier(random_state=42)),
    ("SVM (RBF Kernel)", SVC(kernel='rbf', probability=True, random_state=42)),
    ("Gradient Boosting", GradientBoostingClassifier(random_state=42)),
    ("K-Nearest Neighbors", KNeighborsClassifier(n_neighbors=5)),
    ("XGBoost", XGBClassifier(use_label_encoder=False, eval_metric='mlogloss'),
]

results = []

# Now, we loop through and train models on the transformed data
# Re-split processed data for classification models
X_train_processed_class, X_test_processed_class, y_train_enc_class, y_test_enc_
    X_processed, y_enc, test_size=0.2, random_state=42, stratify=y_enc
)
y_train_class = le.inverse_transform(y_train_enc_class) # Get original labels b
y_test_class = le.inverse_transform(y_test_enc_class)

for name, model in models:
    if name == "XGBoost":

```

```

        model.fit(X_train_processed_class, y_train_enc_class)
        acc = accuracy_score(y_test_enc_class, model.predict(X_test_processed_class))
    else:
        model.fit(X_train_processed_class, y_train_class)
        acc = accuracy_score(y_test_class, model.predict(X_test_processed_class))
    results.append((name, acc))

df_results = pd.DataFrame(results, columns=["Model", "Accuracy"]).sort_values(by="Accuracy", ascending=False)

print("\n==== Model Comparison ===")
print(df_results.to_string(index=False))

# Optional: Print classification report for the best model
best_model_name = df_results.iloc[0]['Model']
best_model = [m for n, m in models if n == best_model_name][0]

print(f"\n==== Classification Report for {best_model_name} ===")
# Need predictions from the best model on the test set corresponding to original labels
if best_model_name == "XGBoost":
    y_pred_class = le.inverse_transform(best_model.predict(X_test_processed_class))
else:
    y_pred_class = best_model.predict(X_test_processed_class)

print(classification_report(y_test_class, y_pred_class))

print("\n--- Classification Model Comparison Complete! ---")

Starting a comprehensive data cleaning and modeling process...
Comprehensive cleaning and feature engineering complete.
Data successfully split into training and testing sets.
Features processed successfully with the specified features.

Shape of X_processed: (6602, 412)
/usr/local/lib/python3.12/dist-packages/xgboost/training.py:183: UserWarning: [1]
Parameters: { "use_label_encoder" } are not used.

bst.update(dtrain, iteration=i, fobj=obj)

==== Model Comparison ===
      Model  Accuracy
Logistic Regression  0.519304
      SVM (RBF Kernel)  0.504164
      Gradient Boosting  0.489780
          XGBoost  0.473883
      Random Forest  0.467827
K-Nearest Neighbors  0.445117

==== Classification Report for Logistic Regression ===
      precision    recall   f1-score   support
AWAY_TEAM         0.50      0.50      0.50       404

```

DRAW	0.31	0.11	0.16	339
HOME_TEAM	0.56	0.77	0.65	578
accuracy			0.52	1321
macro avg	0.46	0.46	0.44	1321
weighted avg	0.48	0.52	0.48	1321

--- Classification Model Comparison Complete! ---

▼ Prepare Data for Regression

Before we train our regression models, we need to make sure our data is in the right shape and format. Think of this like getting your ingredients ready before you start cooking!

For regression, which is about predicting a continuous number (like the difference in scores), our "target" (what we want to predict) needs to be a single number for each match. We already created a 'goal_difference' column earlier, which is perfect for this. It's the home team's score minus the away team's score.

Our "features" (the information we use to make the prediction) are all the other columns we've cleaned and engineered. We used scaling on these features before for classification, which helps models work better, and that scaling is still appropriate for most regression models.

So, in this step, we'll just explicitly define our features (X) and our numerical target (y) for the regression task.

```
# ⚠️ IMPORTANT: Run this cell before training any regression models.
# This block defines the regression features and target variable (goal_difference)
# using only pre-match information to avoid data leakage.

# --- Step 1: Create the target variable for regression ---
# This is the actual goal difference from the match (home - away).
# It's safe to use as a target for regression, but should NOT be used as a feature.
if 'goal_difference' not in cleaned_df.columns:
    cleaned_df['goal_difference'] = cleaned_df['home_score'] - cleaned_df['away_score']

# --- Step 2: Define categorical features ---
# These are identifiers and labels known before the match.
categorical_features_reg = [
    'competition_name', 'competition_code',
    'homeTeam_name', 'homeTeam_tla',
    'awayTeam_name', 'awayTeam_tla'
]
```

```
# --- Step 3: Define numerical features ---
# These are engineered features based on historical data and match metadata.
numerical_features_reg = [
    'expected_goal_diff',                      # Proxy for goal difference based on rol
    'match_hour',                             # Time of match
    'match_day_of_week',                     # Day of week
    'match_month',                           # Month of match
    'is_home_game',                          # Binary indicator for home advantage
    'homeTeam_rolling_avg_goals_scored_5',   # Historical rolling average
    'awayTeam_rolling_avg_goals_conceded_5' # Historical rolling average
]

# --- Step 4: Filter out missing columns ---
# This ensures we only use columns that actually exist in the DataFrame.
numerical_features_reg = [f for f in numerical_features_reg if f in cleaned_df]
categorical_features_reg = [f for f in categorical_features_reg if f in cleaned_df]

# --- Step 5: Combine features and define X and y ---
all_features_reg = categorical_features_reg + numerical_features_reg
X_reg = cleaned_df[all_features_reg].copy()
y_reg = cleaned_df['goal_difference']

# --- Step 6: Preview the data ---
print("✅ Regression features and target defined.")
print(f"Features shape: {X_reg.shape}")
print(f"Target shape: {y_reg.shape}")
print("\n📊 First 5 rows of features:")
display(X_reg.head())
print("\n🎯 First 5 rows of target:")
display(y_reg.head())
```

Most regression models (like Ridge, OLS, Kernel Ridge) cannot handle NaN values in the input features. If you feed them data with missing values, they'll throw errors or silently fail. So before training, you need to either:

- Impute (fill in) missing values with something reasonable (like the mean or median), or
- Drop rows with missing data (which can reduce your sample size).

```
# ⚠ Step: Handle Missing Values in Regression Features
# Many regression models (like LinearRegression, Ridge, etc.) cannot handle NaN
# This step ensures that all numerical features are clean before training.

from sklearn.impute import SimpleImputer

# --- Identify numerical columns in X_reg ---
# These are the columns most likely to contain NaNs and need imputation.
numerical_features_reg = [f for f in X_reg.columns if X_reg[f].dtype in ['float64', 'int64']]

# --- Create a SimpleImputer ---
# We'll use the median strategy, which is robust to outliers and preserves distribution
imputer = SimpleImputer(strategy='median')

# --- Apply imputation to numerical features ---
# This replaces NaNs with the median value of each column.
X_reg[numerical_features_reg] = imputer.fit_transform(X_reg[numerical_features_reg])

# --- Optional: Drop rows with remaining NaNs (e.g., in categorical columns) ---
# This ensures the final dataset is fully clean.
X_reg = X_reg.dropna()
y_reg = y_reg.loc[X_reg.index] # Keep target aligned with filtered features

# --- Confirm cleanup ---
print(f"✅ Missing values handled. Final shape of X_reg: {X_reg.shape}")

# ✅ Missing values handled. Final shape of X_reg: (6602, 7)
```

```

# Check total number of missing values across all columns
total_missing = X_reg.isnull().sum().sum()

# Check how many rows still contain any NaNs
rows_with_nan = X_reg.isnull().any(axis=1).sum()

# Print results
print(f"🔍 Total missing values in X_reg: {total_missing}")
print(f"🔍 Rows with at least one NaN: {rows_with_nan}")

# Final verdict
if total_missing == 0 and rows_with_nan == 0:
    print("✅ All missing values successfully removed or imputed.")
else:
    print("⚠️ Warning: Some missing values still remain. Check preprocessing steps.")

🔍 Total missing values in X_reg: 0
🔍 Rows with at least one NaN: 0
✅ All missing values successfully removed or imputed.

```

▼ Train and Evaluate OLS and Ridge Regression

We're starting with two basic but important types of regression:

- 1. Ordinary Least Squares (OLS):** This is like trying to find the straight line that best fits our data points to predict the goal difference. It's a fundamental method that's easy to understand.
- 2. Ridge Regression:** This is similar to OLS but with a slight twist. It adds a penalty to the model to prevent it from becoming too complex or overly sensitive to individual data points. This can sometimes help the model perform better on new, unseen data.

We'll train both models and see how well they predict the goal difference on our testing data. We'll measure their performance using something called **Mean Squared Error (MSE)**. A lower MSE means the model's predictions are closer to the actual goal differences. We'll also keep track of how long it takes to train each model.

```

from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression, Ridge
from sklearn.metrics import mean_squared_error
import time
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.pipeline import Pipeline

```

```

from sklearn.impute import SimpleImputer # Import SimpleImputer

# Re-split the data specifically for regression features and target
# Use the X_reg and y_reg created in the previous step
X_train_reg, X_test_reg, y_train_reg, y_test_reg = train_test_split(
    X_reg, y_reg, test_size=0.2, random_state=42
)

print("Regression data split into training and testing sets.")
print(f"Training features shape: {X_train_reg.shape}")
print(f"Testing features shape: {X_test_reg.shape}")
print(f"Training target shape: {y_train_reg.shape}")
print(f"Testing target shape: {y_test_reg.shape}")

# Define the preprocessor again, ensuring it matches the features used for reg
# Use the categorical_features_reg and numerical_features_reg lists
preprocessor_reg = ColumnTransformer(
    transformers=[
        ('num', Pipeline([
            ('imputer', SimpleImputer(strategy='constant', fill_value=0)), # Im
            ('scaler', StandardScaler())
        ]), numerical_features_reg),
        ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features_re
    ],
    remainder='passthrough'
)
)

# --- Train and Evaluate OLS ---
print("\nTraining Ordinary Least Squares (OLS) model...")
start_time_lr = time.time()

# Create a pipeline that first preprocesses the data and then applies Linear Re
model_lr = Pipeline(steps=[('preprocessor', preprocessor_reg),
                           ('regressor', LinearRegression())])

model_lr.fit(X_train_reg, y_train_reg)
end_time_lr = time.time()
training_time_lr = end_time_lr - start_time_lr

y_pred_lr = model_lr.predict(X_test_reg)
mse_lr = mean_squared_error(y_test_reg, y_pred_lr)

print(f"OLS Training Time: {training_time_lr:.4f} seconds")
print(f"OLS Validation MSE: {mse_lr:.4f}")

# --- Train and Evaluate Ridge Regression ---
print("\nTraining Ridge Regression model...")

```

```

start_time_ridge = time.time()

# Create a pipeline for Ridge Regression
model_ridge = Pipeline(steps=[('preprocessor', preprocessor_reg),
                             ('regressor', Ridge(alpha=1.0))]) # Using default

model_ridge.fit(X_train_reg, y_train_reg)
end_time_ridge = time.time()
training_time_ridge = end_time_ridge - start_time_ridge

y_pred_ridge = model_ridge.predict(X_test_reg)
mse_ridge = mean_squared_error(y_test_reg, y_pred_ridge)

print(f'Ridge Training Time (alpha=1.0): {training_time_ridge:.4f} seconds')
print(f'Ridge Validation MSE (alpha=1.0): {mse_ridge:.4f}')

# Store results for later comparison
results_regression = {
    "OLS": {"training_time": training_time_lr, "validation_mse": mse_lr},
    "Ridge (alpha=1.0)": {"training_time": training_time_ridge, "validation_mse": mse_ridge}
}

```

Regression data split into training and testing sets.

Training features shape: (5281, 7)

Testing features shape: (1321, 7)

Training target shape: (5281,)

Testing target shape: (1321,)

Training Ordinary Least Squares (OLS) model...

OLS Training Time: 0.0907 seconds

OLS Validation MSE: 2.8760

Training Ridge Regression model...

Ridge Training Time (alpha=1.0): 0.0679 seconds

Ridge Validation MSE (alpha=1.0): 2.8847

✓ Train and Evaluate Kernel Ridge Regression (RBF Kernel)

Next, we'll explore **Kernel Ridge Regression** with a **Radial Basis Function (RBF)** kernel.

This is a more flexible model that can capture non-linear relationships in the data, potentially leading to better predictions than simple linear models like OLS or Ridge.

The RBF kernel uses a parameter called `gamma`, which controls the shape of the kernel. We also have the `alpha` parameter, similar to Ridge Regression, which helps prevent overfitting. We'll try out different combinations of `alpha` and `gamma` to see which ones give us the lowest Mean Squared Error (MSE) on our testing data. This process of trying different parameter values is called **hyperparameter tuning**.

```

# Section: Kernel Ridge Regression (Soccer - Goal Difference)
from sklearn.kernel_ridge import KernelRidge
from sklearn.metrics import mean_squared_error
import time
import pandas as pd # Import pandas to store results
from sklearn.pipeline import Pipeline # Import Pipeline

# Assuming X_train_reg, X_test_reg, y_train_reg, y_test_reg, and preprocessor_rbf

print("\nTraining Kernel Ridge Regression with RBF kernel and tuning hyperparameters")

# Define the hyperparameters to tune
param_grid_rbf = {
    'alpha': [0.1, 1.0, 10.0],
    'gamma': [0.01, 0.1, 1.0]
}

results_kr_rbf = []

# Perform manual grid search
for alpha_val in param_grid_rbf['alpha']:
    for gamma_val in param_grid_rbf['gamma']:
        print(f"\nTraining with alpha={alpha_val}, gamma={gamma_val}")
        start_time_kr_rbf = time.time()

        # Create a pipeline with preprocessing and Kernel Ridge (RBF)
        model_kr_rbf = Pipeline(steps=[('preprocessor', preprocessor_reg),
                                         ('regressor', KernelRidge(kernel='rbf', alpha=alpha_val, gamma=gamma_val))])

        model_kr_rbf.fit(X_train_reg, y_train_reg)
        end_time_kr_rbf = time.time()
        training_time_kr_rbf = end_time_kr_rbf - start_time_kr_rbf

        y_pred_kr_rbf = model_kr_rbf.predict(X_test_reg)
        mse_kr_rbf = mean_squared_error(y_test_reg, y_pred_kr_rbf)

        print(f" Training Time: {training_time_kr_rbf:.4f} seconds")
        print(f" Validation MSE: {mse_kr_rbf:.4f}")

        results_kr_rbf.append({
            "alpha": alpha_val,
            "gamma": gamma_val,
            "training_time": training_time_kr_rbf,
            "validation_mse": mse_kr_rbf
        })

# Store results in a DataFrame
df_kr_rbf = pd.DataFrame(results_kr_rbf)

print("\n--- Kernel Ridge (RBF) Results ---")

```

```

print("---- Kernel Ridge (RBF) results ---",
      print(df_kr_rbf.to_string(index=False))

# Find the best parameters and lowest MSE
best_kr_rbf = df_kr_rbf.loc[df_kr_rbf['validation_mse'].idxmin()]

print(f"\nBest Kernel Ridge (RBF) Model:")
print(f"  Alpha: {best_kr_rbf['alpha']}"))
print(f"  Gamma: {best_kr_rbf['gamma']}"))
print(f"  Lowest Validation MSE: {best_kr_rbf['validation_mse']:.4f}")
print(f"  Training Time for Best Model: {best_kr_rbf['training_time']:.4f} secc

# Add best result to the overall regression results dictionary
results_regression["Kernel Ridge (RBF)"] = {
    "Best Alpha": best_kr_rbf['alpha'],
    "Best Gamma": best_kr_rbf['gamma'],
    "training_time": best_kr_rbf['training_time'],
    "validation_mse": best_kr_rbf['validation_mse']
}

```

Training Kernel Ridge Regression with RBF kernel and tuning hyperparameters...

Training with alpha=0.1, gamma=0.01
 Training Time: 8.6280 seconds
 Validation MSE: 2.9236

Training with alpha=0.1, gamma=0.1
 Training Time: 10.0165 seconds
 Validation MSE: 3.2304

Training with alpha=0.1, gamma=1.0
 Training Time: 10.7877 seconds
 Validation MSE: 3.6308

Training with alpha=1.0, gamma=0.01
 Training Time: 9.4178 seconds
 Validation MSE: 3.0841

Training with alpha=1.0, gamma=0.1
 Training Time: 7.0310 seconds
 Validation MSE: 2.9652

Training with alpha=1.0, gamma=1.0
 Training Time: 8.2232 seconds
 Validation MSE: 3.4778

Training with alpha=10.0, gamma=0.01
 Training Time: 6.7255 seconds
 Validation MSE: 3.2307

Training with alpha=10.0, gamma=0.1
 Training Time: 7.1834 seconds
 Validation MSE: 3.1222

```

Training with alpha=10.0, gamma=1.0
Training Time: 4.3144 seconds
Validation MSE: 3.5837

--- Kernel Ridge (RBF) Results ---
alpha  gamma  training_time  validation_mse
0.1    0.01    8.627978      2.923575
0.1    0.10    10.016532      3.230396
0.1    1.00    10.787679      3.630766
1.0    0.01    9.417840      3.084125
1.0    0.10    7.030954      2.965168
1.0    1.00    8.223183      3.477819
10.0   0.01   6.725468      3.230718
10.0   0.10   7.183378      3.122205
10.0   1.00   4.314396      3.583694

Best Kernel Ridge (RBF) Model:
Alpha: 0.1
Gamma: 0.01
Lowest Validation MSE: 2.9236
Training Time for Best Model: 8.6280 seconds

```

Train and Evaluate Kernel Ridge Regression (Polynomial Kernel)

We'll continue exploring Kernel Ridge Regression, this time using a **Polynomial kernel**. This kernel allows the model to capture non-linear relationships by essentially creating new features that are combinations of the original features raised to different powers.

For the polynomial kernel, we'll tune three hyperparameters:

1. **alpha**: Similar to Ridge Regression, controls the regularization strength.
2. **gamma**: Influences the shape of the polynomial. If `gamma='auto'`, it will use $1/n_features$.
3. **degree**: The degree of the polynomial. A higher degree means more complex interactions between features.

We'll perform another manual grid search over these parameters and track the training time and validation MSE for each combination.

```

from sklearn.kernel_ridge import KernelRidge
from sklearn.metrics import mean_squared_error
import time
import pandas as pd # Import pandas to store results
from sklearn.pipeline import Pipeline # Import Pipeline

```

```

# Assuming X_train_reg, X_test_reg, y_train_reg, y_test_reg, and preprocessor_reg

print("\nTraining Kernel Ridge Regression with Polynomial kernel and tuning hyperparameters")

# Define the hyperparameters to tune for the polynomial kernel
param_grid_poly = {
    'alpha': [0.1, 1.0, 10.0],
    'gamma': [0.01, 0.1, 1.0], # Using specific gamma values instead of 'auto'
    'degree': [2, 3] # Common degrees to test
}

results_kr_poly = []

# Perform manual grid search
for alpha_val in param_grid_poly['alpha']:
    for gamma_val in param_grid_poly['gamma']:
        for degree_val in param_grid_poly['degree']:
            print(f"\nTraining with alpha={alpha_val}, gamma={gamma_val}, degree={degree_val}")
            start_time_kr_poly = time.time()

            # Create a pipeline with preprocessing and Kernel Ridge (Polynomial)
            model_kr_poly = Pipeline(steps=[('preprocessor', preprocessor_reg),
                                             ('regressor', KernelRidge(kernel='poly'))])

            model_kr_poly.fit(X_train_reg, y_train_reg)
            end_time_kr_poly = time.time()
            training_time_kr_poly = end_time_kr_poly - start_time_kr_poly

            y_pred_kr_poly = model_kr_poly.predict(X_test_reg)
            mse_kr_poly = mean_squared_error(y_test_reg, y_pred_kr_poly)

            print(f" Training Time: {training_time_kr_poly:.4f} seconds")
            print(f" Validation MSE: {mse_kr_poly:.4f}")

            results_kr_poly.append({
                "alpha": alpha_val,
                "gamma": gamma_val,
                "degree": degree_val,
                "training_time": training_time_kr_poly,
                "validation_mse": mse_kr_poly
            })

# Store results in a DataFrame
df_kr_poly = pd.DataFrame(results_kr_poly)

print("\n--- Kernel Ridge (Polynomial) Results ---")
print(df_kr_poly.to_string(index=False))

# Find the best parameters and lowest MSE
best_kr_poly = df_kr_poly.loc[df_kr_poly['validation_mse'].idxmin()]

```

```

print(f"\nBest Kernel Ridge (Polynomial) Model:")
print(f"  Alpha: {best_kr_poly['alpha']}")
print(f"  Gamma: {best_kr_poly['gamma']}")
print(f"  Degree: {best_kr_poly['degree']}")
print(f"  Lowest Validation MSE: {best_kr_poly['validation_mse']:.4f}")
print(f"  Training Time for Best Model: {best_kr_poly['training_time']:.4f} sec

# Add best result to the overall regression results dictionary
results_regression["Kernel Ridge (Poly)"] = {
    "Best Alpha": best_kr_poly['alpha'],
    "Best Gamma": best_kr_poly['gamma'],
    "Best Degree": best_kr_poly['degree'],
    "training_time": best_kr_poly['training_time'],
    "validation_mse": best_kr_poly['validation_mse']
}

```

Training Kernel Ridge Regression with Polynomial kernel and tuning hyperparameters

Training with alpha=0.1, gamma=0.01, degree=2
 Training Time: 5.2552 seconds
 Validation MSE: 2.9210

Training with alpha=0.1, gamma=0.01, degree=3
 Training Time: 4.6995 seconds
 Validation MSE: 2.9102

Training with alpha=0.1, gamma=0.1, degree=2
 Training Time: 4.0055 seconds
 Validation MSE: 3.2325

Training with alpha=0.1, gamma=0.1, degree=3
 Training Time: 5.3255 seconds
 Validation MSE: 3.8270

Training with alpha=0.1, gamma=1.0, degree=2
 Training Time: 4.0307 seconds
 Validation MSE: 4.9253

Training with alpha=0.1, gamma=1.0, degree=3
 Training Time: 6.9433 seconds
 Validation MSE: 13.9476

Training with alpha=1.0, gamma=0.01, degree=2
 Training Time: 3.9923 seconds
 Validation MSE: 3.0795

Training with alpha=1.0, gamma=0.01, degree=3
 Training Time: 4.7022 seconds
 Validation MSE: 3.0385

Training with alpha=1.0, gamma=0.1, degree=2
 Training Time: 4.7805 seconds
 Validation MSE: 3.0000

```

VALIDATION MSE: 2.9908

Training with alpha=1.0, gamma=0.1, degree=3
Training Time: 5.3597 seconds
Validation MSE: 3.0548

Training with alpha=1.0, gamma=1.0, degree=2
Training Time: 5.5910 seconds
Validation MSE: 4.2008

Training with alpha=1.0, gamma=1.0, degree=3
Training Time: 4.8429 seconds
Validation MSE: 7.1638

Training with alpha=10.0, gamma=0.01, degree=2
Training Time: 4.0279 seconds
Validation MSE: 3.2302

Training with alpha=10.0, gamma=0.01, degree=3
Training Time: 5.7238 seconds
Validation MSE: 3.2112

```

```

from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression, Ridge
from sklearn.metrics import mean_squared_error
import time
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer # Import SimpleImputer

# Re-split the data specifically for regression features and target
# Use the X_reg and y_reg created in the previous step
X_train_reg, X_test_reg, y_train_reg, y_test_reg = train_test_split(
    X_reg, y_reg, test_size=0.2, random_state=42
)

print("Regression data split into training and testing sets.")
print(f"Training features shape: {X_train_reg.shape}")
print(f"Testing features shape: {X_test_reg.shape}")
print(f"Training target shape: {y_train_reg.shape}")
print(f"Testing target shape: {y_test_reg.shape}")

# Define the preprocessor again, ensuring it matches the features used for regre
# Use the categorical_features_reg and numerical_features_reg lists
preprocessor_reg = ColumnTransformer(
    transformers=[
        ('num', Pipeline([
            ('imputer', SimpleImputer(strategy='constant', fill_value=0)), # Im
            ('scaler', StandardScaler())
        ]), numerical_features_reg)
    ]
)

```

```

        ],
        ('cat', OneHotEncoder(handle_unknown='ignore')), categorical_features_reg
    ],
    remainder='passthrough'
)

# --- Train and Evaluate OLS ---
print("\nTraining Ordinary Least Squares (OLS) model...")
start_time_lr = time.time()

# Create a pipeline that first preprocesses the data and then applies Linear Regression
model_lr = Pipeline(steps=[('preprocessor', preprocessor_reg),
                           ('regressor', LinearRegression())])

model_lr.fit(X_train_reg, y_train_reg)
end_time_lr = time.time()
training_time_lr = end_time_lr - start_time_lr

y_pred_lr = model_lr.predict(X_test_reg)
mse_lr = mean_squared_error(y_test_reg, y_pred_lr)

print(f"OLS Training Time: {training_time_lr:.4f} seconds")
print(f"OLS Validation MSE: {mse_lr:.4f}")

# --- Train and Evaluate Ridge Regression ---
print("\nTraining Ridge Regression model...")
start_time_ridge = time.time()

# Create a pipeline for Ridge Regression
model_ridge = Pipeline(steps=[('preprocessor', preprocessor_reg),
                               ('regressor', Ridge(alpha=1.0))]) # Using default

model_ridge.fit(X_train_reg, y_train_reg)
end_time_ridge = time.time()
training_time_ridge = end_time_ridge - start_time_ridge

y_pred_ridge = model_ridge.predict(X_test_reg)
mse_ridge = mean_squared_error(y_test_reg, y_pred_ridge)

print(f"Ridge Training Time (alpha=1.0): {training_time_ridge:.4f} seconds")
print(f"Ridge Validation MSE (alpha=1.0): {mse_ridge:.4f}")

# Store results for later comparison
results_regression = {
    "OLS": {"training_time": training_time_lr, "validation_mse": mse_lr},
    "Ridge (alpha=1.0)": {"training_time": training_time_ridge, "validation_mse": mse_ridge}
}

```

Regression data split into training and testing sets.

```

Training features shape: (5281, 7)
Testing features shape: (1321, 7)
Training target shape: (5281,)
Testing target shape: (1321,)

Training Ordinary Least Squares (OLS) model...
OLS Training Time: 0.0680 seconds
OLS Validation MSE: 2.8760

Training Ridge Regression model...
Ridge Training Time (alpha=1.0): 0.0451 seconds
Ridge Validation MSE (alpha=1.0): 2.8847

```

```

# Assuming 'cleaned_df' is the DataFrame with engineered features from previous

# Define the features (X) - these are the columns we'll use to predict the goal
# We'll exclude the actual goal difference and individual scores from the features
# as those would leak information about the outcome. We'll include the engineered
# like 'expected_goal_diff', date/time features, and team/competition identifiers

# Let's select the relevant features that would be available *before* the match
# We'll use the categorical features and the engineered numerical features like
categorical_features_reg = [
    'competition_name', 'competition_code',
    'homeTeam_name', 'homeTeam_tla',
    'awayTeam_name', 'awayTeam_tla'
]

numerical_features_reg = [
    'expected_goal_diff', # This is our pre-match proxy for goal difference
    'match_hour',
    'match_day_of_week',
    'match_month',
    'is_home_game', # This is our binary home advantage indicator
    'homeTeam_rolling_avg_goals_scored_5', # Rolling averages from previous code
    'awayTeam_rolling_avg_goals_conceded_5' # Rolling averages from previous code
]

# Filter out any features that might not exist in the DataFrame based on previous
numerical_features_reg = [f for f in numerical_features_reg if f in cleaned_df]
categorical_features_reg = [f for f in categorical_features_reg if f in cleaned_df]

# Combine the feature lists
all_features_reg = categorical_features_reg + numerical_features_reg

# Create the features DataFrame
X_reg = cleaned_df[all_features_reg].copy()

# Define the target variable (y) - this is what we want to predict.
y_reg = cleaned_df['goal_difference']

```

```
print("Features (X_reg) and target (y_reg) for regression defined.")
print(f"Features shape: {X_reg.shape}")
print(f"Target shape: {y_reg.shape}")
print("\nFirst 5 rows of features:")
display(X_reg.head())
print("\nFirst 5 rows of target:")
display(y_reg.head())
```

Train and Evaluate Kernel Ridge Regression (Polynomial Kernel)

We'll continue exploring Kernel Ridge Regression, this time using a **Polynomial kernel**. This kernel allows the model to capture non-linear relationships by essentially creating new features that are combinations of the original features raised to different powers.

For the polynomial kernel, we'll tune three hyperparameters:

1. **alpha** : Similar to Ridge Regression, controls the regularization strength.
2. **gamma** : Influences the shape of the polynomial. If `gamma='auto'`, it will use $1/n_features$.
3. **degree** : The degree of the polynomial. A higher degree means more complex interactions between features.

We'll perform another manual grid search over these parameters and track the training time and validation MSE for each combination.

```

from sklearn.kernel_ridge import KernelRidge
from sklearn.metrics import mean_squared_error
import time
import pandas as pd # Import pandas to store results
from sklearn.pipeline import Pipeline # Import Pipeline

# Assuming X_train_reg, X_test_reg, y_train_reg, y_test_reg, and preprocessor_r

print("\nTraining Kernel Ridge Regression with Polynomial kernel and tuning hyp")

# Define the hyperparameters to tune for the polynomial kernel
param_grid_poly = {
    'alpha': [0.1, 1.0, 10.0],
    'gamma': [0.01, 0.1, 1.0], # Using specific gamma values instead of 'auto'
    'degree': [2, 3] # Common degrees to test
}

results_kr_poly = []

# Perform manual grid search
for alpha_val in param_grid_poly['alpha']:
    for gamma_val in param_grid_poly['gamma']:
        for degree_val in param_grid_poly['degree']:
            print(f"\nTraining with alpha={alpha_val}, gamma={gamma_val}, degree={degree_val}")
            start_time_kr_poly = time.time()

            # Create a pipeline with preprocessing and Kernel Ridge (Polynomial)
            model_kr_poly = Pipeline(steps=[('preprocessor', preprocessor_reg),
                                             ('regressor', KernelRidge(kernel='poly'))])

            model_kr_poly.fit(X_train_reg, y_train_reg)
            end_time_kr_poly = time.time()
            training_time_kr_poly = end_time_kr_poly - start_time_kr_poly

            y_pred_kr_poly = model_kr_poly.predict(X_test_reg)
            mse_kr_poly = mean_squared_error(y_test_reg, y_pred_kr_poly)

            print(f" Training Time: {training_time_kr_poly:.4f} seconds")
            print(f" Validation MSE: {mse_kr_poly:.4f}\n")

```

```

print(f"\nTraining with alpha={alpha}, gamma={gamma}, degree={degree}\n")
print(f"Training Time: {training_time:.4f} seconds\n")
print(f"Validation MSE: {mse:.4f}\n")

results_kr_poly.append({
    "alpha": alpha_val,
    "gamma": gamma_val,
    "degree": degree_val,
    "training_time": training_time_kr_poly,
    "validation_mse": mse_kr_poly
})

# Store results in a DataFrame
df_kr_poly = pd.DataFrame(results_kr_poly)

print("\n--- Kernel Ridge (Polynomial) Results ---")
print(df_kr_poly.to_string(index=False))

# Find the best parameters and lowest MSE
best_kr_poly = df_kr_poly.loc[df_kr_poly['validation_mse'].idxmin()]

print(f"\nBest Kernel Ridge (Polynomial) Model:")
print(f"  Alpha: {best_kr_poly['alpha']} ")
print(f"  Gamma: {best_kr_poly['gamma']} ")
print(f"  Degree: {best_kr_poly['degree']} ")
print(f"  Lowest Validation MSE: {best_kr_poly['validation_mse']:.4f}")
print(f"  Training Time for Best Model: {best_kr_poly['training_time']:.4f} sec")

# Add best result to the overall regression results dictionary
results_regression["Kernel Ridge (Poly)"] = {
    "Best Alpha": best_kr_poly['alpha'],
    "Best Gamma": best_kr_poly['gamma'],
    "Best Degree": best_kr_poly['degree'],
    "training_time": best_kr_poly['training_time'],
    "validation_mse": best_kr_poly['validation_mse']
}

```

Training Kernel Ridge Regression with Polynomial kernel and tuning hyperparameters

Training with alpha=0.1, gamma=0.01, degree=2

Training Time: 5.8494 seconds

Validation MSE: 2.9210

Training with alpha=0.1, gamma=0.01, degree=3

Training Time: 5.9410 seconds

Validation MSE: 2.9102

Training with alpha=0.1, gamma=0.1, degree=2

Training Time: 4.1825 seconds

Validation MSE: 3.2325

Training with alpha=0.1, gamma=0.1, degree=3

Training Time: 5.2628 seconds

Validation MSE: 3.8270

VALIDATION MSE: 3.0270

Training with alpha=0.1, gamma=1.0, degree=2

Training Time: 3.9805 seconds

Validation MSE: 4.9253

Training with alpha=0.1, gamma=1.0, degree=3

Training Time: 6.6240 seconds

Validation MSE: 13.9476

Training with alpha=1.0, gamma=0.01, degree=2

Training Time: 4.0199 seconds

Validation MSE: 3.0795

Training with alpha=1.0, gamma=0.01, degree=3

Training Time: 4.6632 seconds

Validation MSE: 3.0385

Training with alpha=1.0, gamma=0.1, degree=2

Training Time: 4.7230 seconds

Validation MSE: 2.9383

Training with alpha=1.0, gamma=0.1, degree=3

Training Time: 4.6796 seconds

Validation MSE: 3.0548

Training with alpha=1.0, gamma=1.0, degree=2

Training Time: 5.0051 seconds

Validation MSE: 4.2008

Training with alpha=1.0, gamma=1.0, degree=3

Training Time: 4.9878 seconds

Validation MSE: 7.1638

Training with alpha=10.0, gamma=0.01, degree=2

Training Time: 4.0478 seconds

Validation MSE: 3.2302

Training with alpha=10.0, gamma=0.01, degree=3

Training Time: 6.0337 seconds

Validation MSE: 3.2112

```
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression, Ridge
from sklearn.metrics import mean_squared_error
import time
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer # Import SimpleImputer
```

```
# Re-split the data specifically for regression features and target
# Use the X_reg and y_reg created in the previous step
.....
```

```

x_train_reg, x_test_reg, y_train_reg, y_test_reg = train_test_split(
    X_reg, y_reg, test_size=0.2, random_state=42
)

print("Regression data split into training and testing sets.")
print(f"Training features shape: {X_train_reg.shape}")
print(f"Testing features shape: {X_test_reg.shape}")
print(f"Training target shape: {y_train_reg.shape}")
print(f"Testing target shape: {y_test_reg.shape}")

# Define the preprocessor again, ensuring it matches the features used for regr
# Use the categorical_features_reg and numerical_features_reg lists
preprocessor_reg = ColumnTransformer(
    transformers=[
        ('num', Pipeline([
            ('imputer', SimpleImputer(strategy='constant', fill_value=0)), # Im
            ('scaler', StandardScaler())
        ]), numerical_features_reg),
        ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features_re
    ],
    remainder='passthrough'
)

# --- Train and Evaluate OLS ---
print("\nTraining Ordinary Least Squares (OLS) model...")
start_time_lr = time.time()

# Create a pipeline that first preprocesses the data and then applies Linear Re
model_lr = Pipeline(steps=[('preprocessor', preprocessor_reg),
                           ('regressor', LinearRegression())])

model_lr.fit(X_train_reg, y_train_reg)
end_time_lr = time.time()
training_time_lr = end_time_lr - start_time_lr

y_pred_lr = model_lr.predict(X_test_reg)
mse_lr = mean_squared_error(y_test_reg, y_pred_lr)

print(f"OLS Training Time: {training_time_lr:.4f} seconds")
print(f"OLS Validation MSE: {mse_lr:.4f}")

# --- Train and Evaluate Ridge Regression ---
print("\nTraining Ridge Regression model...")
start_time_ridge = time.time()

# Create a pipeline for Ridge Regression
model_ridge = Pipeline(steps=[('preprocessor', preprocessor_reg),
                             ('regressor', Ridge(alpha=1.0))]) # Using default

```

```

model_ridge.fit(X_train_reg, y_train_reg)
end_time_ridge = time.time()
training_time_ridge = end_time_ridge - start_time_ridge

y_pred_ridge = model_ridge.predict(X_test_reg)
mse_ridge = mean_squared_error(y_test_reg, y_pred_ridge)

print(f"Ridge Training Time (alpha=1.0): {training_time_ridge:.4f} seconds")
print(f"Ridge Validation MSE (alpha=1.0): {mse_ridge:.4f}")

# Store results for later comparison
results_regression = {
    "OLS": {"training_time": training_time_lr, "validation_mse": mse_lr},
    "Ridge (alpha=1.0)": {"training_time": training_time_ridge, "validation_mse": mse_ridge}
}

```

Regression data split into training and testing sets.

Training features shape: (5281, 7)

Testing features shape: (1321, 7)

Training target shape: (5281,)

Testing target shape: (1321,)

Training Ordinary Least Squares (OLS) model...

OLS Training Time: 0.0662 seconds

OLS Validation MSE: 2.8759

Training Ridge Regression model...

Ridge Training Time (alpha=1.0): 0.0415 seconds

Ridge Validation MSE (alpha=1.0): 2.8848

```

# Assuming 'cleaned_df' is the DataFrame with engineered features from previous

# Define the features (X) - these are the columns we'll use to predict the goal
# We'll exclude the actual goal difference and individual scores from the features
# as those would leak information about the outcome. We'll include the engineered
# features like 'expected_goal_diff', date/time features, and team/competition identifiers

# Let's select the relevant features that would be available *before* the match
# We'll use the categorical features and the engineered numerical features like
categorical_features_reg = [
    'competition_name', 'competition_code',
    'homeTeam_name', 'homeTeam_tla',
    'awayTeam_name', 'awayTeam_tla'
]

numerical_features_reg = [
    'expected_goal_diff', # This is our pre-match proxy for goal difference
    'match_hour',
    'match_day_of_week',
    'match_month',

```

```
'is_home_game', # This is our binary home advantage indicator
'homeTeam_rolling_avg_goals_scored_5', # Rolling averages from previous cod
'awayTeam_rolling_avg_goals_conceded_5' # Rolling averages from previous cc
]

# Filter out any features that might not exist in the DataFrame based on previous
numerical_features_reg = [f for f in numerical_features_reg if f in cleaned_df]
categorical_features_reg = [f for f in categorical_features_reg if f in cleaned_df]

# Combine the feature lists
all_features_reg = categorical_features_reg + numerical_features_reg

# Create the features DataFrame
X_reg = cleaned_df[all_features_reg].copy()

# Define the target variable (y) - this is what we want to predict.
y_reg = cleaned_df['goal_difference']

print("Features (X_reg) and target (y_reg) for regression defined.")
print(f"Features shape: {X_reg.shape}")
print(f"Target shape: {y_reg.shape}")
print("\nFirst 5 rows of features:")
display(X_reg.head())
print("\nFirst 5 rows of target:")
display(y_reg.head())
```

```

import requests
import pandas as pd
import time

# YOUR API KEY GOES HERE
api_key = "e0d3c18b3e2c496e9a79bf92ccbbe53f"
BASE_URL = "https://api.football-data.org/v4/"

headers = {
    'X-Auth-Token': api_key,
    'X-Unfold-Goals': 'true'
}

# The leagues to include, based on your free plan
competition_codes = ['BL1', 'DED', 'BSA', 'PD', 'FL1', 'ELC', 'PPL', 'SA', 'PL'
seasons = [2023, 2024]

all_matches = []
unique_match_ids = set()

print("Starting data collection for multiple leagues and seasons...")
for competition_code in competition_codes:
    for season in seasons:
        print(f"\nFetching matches for {competition_code} in {season}/{season+1}")
        url = f"{BASE_URL}competitions/{competition_code}/matches?season={seasc

        try:
            response = requests.get(url, headers=headers)
            response.raise_for_status()
            data = response.json()

            if 'matches' in data:
                for match in data['matches']:
                    if match['id'] not in unique_match_ids:
                        all_matches.append(match)
                        unique_match_ids.add(match['id'])
                print(f" Successfully fetched {len(data['matches'])} matches.")
            else:
                print(f" No matches found for {competition_code} in that seasc

        except requests.exceptions.RequestException as e:
            print(f" An error occurred: {e}")

```

```
    time.sleep(15) # Respect the API rate limit

# --- Create a Final DataFrame ---
print("\nCreating a single DataFrame from the collected data...")
if all_matches:
    final_df = pd.DataFrame(all_matches)

    print(f"\nTotal unique matches in the final DataFrame: {len(final_df)}")

    print("\nDataFrame created. Here are the first 5 rows:")
    print(final_df.head())

    print("\nDataFrame Columns:")
    print(final_df.columns)
else:
    print("No match data was collected.")
```

Starting data collection for multiple leagues and seasons...

Fetching matches for BL1 in 2023/2024 season...
Successfully fetched 306 matches.

Fetching matches for BL1 in 2024/2025 season...
Successfully fetched 306 matches.

Fetching matches for DED in 2023/2024 season...
Successfully fetched 306 matches.

Fetching matches for DED in 2024/2025 season...
Successfully fetched 306 matches.

Fetching matches for BSA in 2023/2024 season...
Successfully fetched 380 matches.

Fetching matches for BSA in 2024/2025 season...
Successfully fetched 380 matches.

Fetching matches for PD in 2023/2024 season...
Successfully fetched 380 matches.

Fetching matches for PD in 2024/2025 season...
Successfully fetched 380 matches.

Fetching matches for FL1 in 2023/2024 season...
Successfully fetched 306 matches.

Fetching matches for FL1 in 2024/2025 season...
Successfully fetched 306 matches.

Fetching matches for ELC in 2023/2024 season...
Successfully fetched 557 matches.

Fetching matches for ELC in 2024/2025 season...
Successfully fetched 557 matches.

```
Fetching matches for PPL in 2023/2024 season...
Successfully fetched 306 matches.
```

```
Fetching matches for PPL in 2024/2025 season...
Successfully fetched 306 matches.
```

```
Fetching matches for SA in 2023/2024 season...
Successfully fetched 380 matches.
```

```
Fetching matches for SA in 2024/2025 season...
Successfully fetched 380 matches.
```

```
Fetching matches for PL in 2023/2024 season...
Successfully fetched 380 matches.
```

```
Fetching matches for PL in 2024/2025 season...
Successfully fetched 380 matches.
```

```
Creating a single DataFrame from the collected data...
```

```
# --- Step 1: Flattening the Data ---
# Extract key information from the nested columns.

# Normalize home and away team data
home_teams = pd.json_normalize(final_df['homeTeam']).add_prefix('homeTeam_')
away_teams = pd.json_normalize(final_df['awayTeam']).add_prefix('awayTeam_')

# Normalize score data
scores = pd.json_normalize(final_df['score'])

# Normalize competition data to get the league name
competitions = pd.json_normalize(final_df['competition']).add_prefix('competiti')

# Combine all the flattened data into a new DataFrame
cleaned_df = pd.concat([
    final_df[['id', 'utcDate', 'status']], # Keep essential columns
    competitions[['competition_name', 'competition_code']],
    home_teams[['homeTeam_name', 'homeTeam_tla']],
    away_teams[['awayTeam_name', 'awayTeam_tla']],
    scores[['winner', 'fullTime.home', 'fullTime.away', 'halfTime.home', 'halfT'],
], axis=1)

# --- Step 2: Handle Missing Values ---
# Fill the few missing halftime scores with 0
cleaned_df['halfTime.home'] = cleaned_df['halfTime.home'].fillna(0).astype(int)
cleaned_df['halfTime.away'] = cleaned_df['halfTime.away'].fillna(0).astype(int)

# Drop columns with too many missing values
columns_to_drop = ['regularTime.home', 'regularTime.away', 'extraTime.home', 'e
cleaned_df.drop(columns=columns_to_drop, inplace=True, errors='ignore')
```

```

# --- Step 3: Rename Columns for Clarity ---
cleaned_df.rename(columns={
    'fullTime.home': 'home_score',
    'fullTime.away': 'away_score',
    'halfTime.home': 'halfTime_home_score',
    'halfTime.away': 'halfTime_away_score'
}, inplace=True)

print("Comprehensive cleaning complete.\n")
print(f"DataFrame now has {len(cleaned_df.columns)} columns and {len(cleaned_df)} rows")
print("The first 5 rows of the new, clean DataFrame:")
print(cleaned_df.head())
print("\nNew DataFrame Columns:")
print(cleaned_df.columns)

```

Comprehensive cleaning complete.

DataFrame now has 14 columns and 6602 rows.

The first 5 rows of the new, clean DataFrame:

	<code>id</code>	<code>utcDate</code>	<code>status</code>	<code>competition_name</code>	<code>competition_code</code>	
0	441789	2023-08-18T18:30:00Z	FINISHED	Bundesliga	BL1	
1	441792	2023-08-19T13:30:00Z	FINISHED	Bundesliga	BL1	
2	441794	2023-08-19T13:30:00Z	FINISHED	Bundesliga	BL1	
3	441795	2023-08-19T13:30:00Z	FINISHED	Bundesliga	BL1	
4	441796	2023-08-19T13:30:00Z	FINISHED	Bundesliga	BL1	

	<code>homeTeam_name</code>	<code>homeTeam_tla</code>	<code>awayTeam_name</code>	<code>awayTeam_tla</code>	
0	SV Werder Bremen	SVW	FC Bayern München	FCB	
1	Bayer 04 Leverkusen	B04	RB Leipzig	RBL	
2	VfL Wolfsburg	WOB	1. FC Heidenheim 1846	HEI	
3	TSG 1899 Hoffenheim	TSG	SC Freiburg	SCF	
4	FC Augsburg	FCA	Borussia Mönchengladbach	BMG	

	<code>winner</code>	<code>home_score</code>	<code>away_score</code>	<code>halfTime_home_score</code>	<code>halfTime_away_score</code>	
0	AWAY_TEAM	0	4	0	1	
1	HOME_TEAM	3	2	2	1	
2	HOME_TEAM	2	0	2	0	
3	AWAY_TEAM	1	2	0	2	
4	DRAW	4	4	3	3	

New DataFrame Columns:

```

Index(['id', 'utcDate', 'status', 'competition_name', 'competition_code',
       'homeTeam_name', 'homeTeam_tla', 'awayTeam_name', 'awayTeam_tla',
       'winner', 'home_score', 'away_score', 'halfTime_home_score',
       'halfTime_away_score'],
      dtype='object')

```

```
# --- Create safe pre-match proxy for goal_difference ---
```

```
# Calculate per-match goal differences (post-match, but we'll use them only for
cleaned_df['home_goal_diff'] = cleaned_df['home_score'] - cleaned_df['away_score']
cleaned_df['away_goal_diff'] = cleaned_df['away_score'] - cleaned_df['home_score']
```

```

# Sort by date so rolling windows are correct
cleaned_df = cleaned_df.sort_values(by='utcDate')

# Rolling averages for last 5 matches (excluding current match)
cleaned_df['home_avg_goal_diff_last5'] = (
    cleaned_df.groupby('homeTeam_name')['home_goal_diff']
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)

cleaned_df['away_avg_goal_diff_last5'] = (
    cleaned_df.groupby('awayTeam_name')['away_goal_diff']
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)

# Combine into a single "expected goal diff" feature
cleaned_df['expected_goal_diff'] = (
    cleaned_df['home_avg_goal_diff_last5'] - cleaned_df['away_avg_goal_diff_last5']
)

# Convert the 'utcDate' column to datetime objects
cleaned_df['utcDate'] = pd.to_datetime(cleaned_df['utcDate'])

# Extract date and time components
cleaned_df['match_hour'] = cleaned_df['utcDate'].dt.hour
cleaned_df['match_day_of_week'] = cleaned_df['utcDate'].dt.dayofweek # Monday=0
cleaned_df['match_month'] = cleaned_df['utcDate'].dt.month

# Create a binary feature for home advantage (1 if home game, 0 otherwise)
cleaned_df['is_home_game'] = 1 # In this dataset structure, each row is a match

# Add rolling averages for goals scored and conceded for home and away teams
# (Assuming these were intended to be created from previous steps based on the
cleaned_df['homeTeam_rolling_avg_goals_scored_5'] = (
    cleaned_df.groupby('homeTeam_name')['home_score']
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)

cleaned_df['awayTeam_rolling_avg_goals_conceded_5'] = (
    cleaned_df.groupby('awayTeam_name')['home_score'] # Away team conceded is home score
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)

print("\nEngineered features added to cleaned_df.")
display(cleaned_df.head())

```

```

# --- Feature Engineering Block ---
cleaned_df['goal_difference'] = cleaned_df['home_score'] - cleaned_df['away_score']
cleaned_df['total_goals'] = cleaned_df['home_score'] + cleaned_df['away_score']

cleaned_df['home_goal_diff'] = cleaned_df['home_score'] - cleaned_df['away_score']
cleaned_df['away_goal_diff'] = cleaned_df['away_score'] - cleaned_df['home_score']

cleaned_df = cleaned_df.sort_values(by='utcDate')

cleaned_df['home_avg_goal_diff_last5'] = (
    cleaned_df.groupby('homeTeam_name')['home_goal_diff']
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)

cleaned_df['away_avg_goal_diff_last5'] = (
    cleaned_df.groupby('awayTeam_name')['away_goal_diff']
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)

cleaned_df['expected_goal_diff'] = (
    cleaned_df['home_avg_goal_diff_last5'] - cleaned_df['away_avg_goal_diff_last5']
)

cleaned_df['match_hour'] = cleaned_df['utcDate'].dt.hour
cleaned_df['match_day_of_week'] = cleaned_df['utcDate'].dt.dayofweek
cleaned_df['match_month'] = cleaned_df['utcDate'].dt.month

```

```

# Define the features (X) - these are the columns we'll use to predict the goal
# We'll exclude the actual goal difference and individual scores from the features
# as those would leak information about the outcome. We'll include the engineered
# features like 'expected_goal_diff', date/time features, and team/competition identifiers

```

```
# Let's select the relevant features that would be available *before* the match
# We'll use the categorical features and the engineered numerical features like
categorical_features_reg = [
    'competition_name', 'competition_code',
    'homeTeam_name', 'homeTeam_tla',
    'awayTeam_name', 'awayTeam_tla'
]

numerical_features_reg = [
    'expected_goal_diff', # This is our pre-match proxy for goal difference
    'match_hour',
    'match_day_of_week',
    'match_month',
    'is_home_game', # This is our binary home advantage indicator
    'homeTeam_rolling_avg_goals_scored_5', # Rolling averages from previous cod
    'awayTeam_rolling_avg_goals_conceded_5' # Rolling averages from previous cc
]

# Filter out any features that might not exist in the DataFrame based on previous
numerical_features_reg = [f for f in numerical_features_reg if f in cleaned_df.
categorical_features_reg = [f for f in categorical_features_reg if f in cleaned

# Combine the feature lists
all_features_reg = categorical_features_reg + numerical_features_reg

# Create the features DataFrame
X_reg = cleaned_df[all_features_reg].copy()

# Define the target variable (y) - this is what we want to predict.
y_reg = cleaned_df['goal_difference']

print("Features (X_reg) and target (y_reg) for regression defined.")
print(f"Features shape: {X_reg.shape}")
print(f"Target shape: {y_reg.shape}")
print("\nFirst 5 rows of features:")
display(X_reg.head())
print("\nFirst 5 rows of target:")
display(y_reg.head())
```

```
# --- Create safe pre-match proxy for goal_difference ---
# Calculate per-match goal differences (post-match, but we'll use them only for
cleaned_df['home_goal_diff'] = cleaned_df['home_score'] - cleaned_df['away_score']
cleaned_df['away_goal_diff'] = cleaned_df['away_score'] - cleaned_df['home_score']

# Sort by date so rolling windows are correct
cleaned_df = cleaned_df.sort_values(by='utcDate')

# Rolling averages for last 5 matches (excluding current match)
cleaned_df['home_avg_goal_diff_last5'] = (
    cleaned_df.groupby('homeTeam_name')['home_goal_diff']
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)

cleaned_df['away_avg_goal_diff_last5'] = (
    cleaned_df.groupby('awayTeam_name')['away_goal_diff']
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)

# Combine into a single "expected goal diff" feature
cleaned_df['expected_goal_diff'] = (
    cleaned_df['home_avg_goal_diff_last5'] - cleaned_df['away_avg_goal_diff_last5']
)

# Convert the 'utcDate' column to datetime objects
```

```

# Convert the utcDate column to datetime objects
cleaned_df['utcDate'] = pd.to_datetime(cleaned_df['utcDate'])

# Extract date and time components
cleaned_df['match_hour'] = cleaned_df['utcDate'].dt.hour
cleaned_df['match_day_of_week'] = cleaned_df['utcDate'].dt.dayofweek # Monday=0
cleaned_df['match_month'] = cleaned_df['utcDate'].dt.month

# Create a binary feature for home advantage (1 if home game, 0 otherwise)
cleaned_df['is_home_game'] = 1 # In this dataset structure, each row is a match

# Add rolling averages for goals scored and conceded for home and away teams
# (Assuming these were intended to be created from previous steps based on the
cleaned_df['homeTeam_rolling_avg_goals_scored_5'] = (
    cleaned_df.groupby('homeTeam_name')['home_score']
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)

cleaned_df['awayTeam_rolling_avg_goals_conceded_5'] = (
    cleaned_df.groupby('awayTeam_name')['home_score'] # Away team conceded is home score
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)

print("\nEngineered features added to cleaned_df.")
display(cleaned_df.head())

```

```

# Define the features (X) - these are the columns we'll use to predict the goal
# We'll exclude the actual goal difference and individual scores from the features
# as those would leak information about the outcome. We'll include the engineered
# features like 'expected goal diff' and date/time features and team/competition identifiers

```

```
# TIME EXPECTED_EVAL_UNIT , DATE/TIME FEATURES, AND TEAM/COMPETITION IDENTIFIERS
```

```
# Let's select the relevant features that would be available *before* the match
# We'll use the categorical features and the engineered numerical features like
categorical_features_reg = [
    'competition_name', 'competition_code',
    'homeTeam_name', 'homeTeam_tla',
    'awayTeam_name', 'awayTeam_tla'
]

numerical_features_reg = [
    'expected_goal_diff', # This is our pre-match proxy for goal difference
    'match_hour',
    'match_day_of_week',
    'match_month',
    'is_home_game', # This is our binary home advantage indicator
    'homeTeam_rolling_avg_goals_scored_5', # Rolling averages from previous cod
    'awayTeam_rolling_avg_goals_conceded_5' # Rolling averages from previous cc
]

# Filter out any features that might not exist in the DataFrame based on previous code
numerical_features_reg = [f for f in numerical_features_reg if f in cleaned_df]
categorical_features_reg = [f for f in categorical_features_reg if f in cleaned_df]

# Combine the feature lists
all_features_reg = categorical_features_reg + numerical_features_reg

# Create the features DataFrame
X_reg = cleaned_df[all_features_reg].copy()

# Define the target variable (y) - this is what we want to predict.
y_reg = cleaned_df['goal_difference']

print("Features (X_reg) and target (y_reg) for regression defined.")
print(f"Features shape: {X_reg.shape}")
print(f"Target shape: {y_reg.shape}")
print("\nFirst 5 rows of features:")
display(X_reg.head())
print("\nFirst 5 rows of target:")
display(y_reg.head())
```

Train and Evaluate Kernel Ridge Regression (Polynomial Kernel)

We'll continue exploring Kernel Ridge Regression, this time using a **Polynomial kernel**. This kernel allows the model to capture non-linear relationships by essentially creating new features that are combinations of the original features raised to different powers.

For the polynomial kernel, we'll tune three hyperparameters:

1. `alpha` : Similar to Ridge Regression, controls the regularization strength.
2. `gamma` : Influences the shape of the polynomial. If `gamma='auto'`, it will use $1/n_features$.
3. `degree` : The degree of the polynomial. A higher degree means more complex interactions between features.

We'll perform another manual grid search over these parameters and track the training time and validation MSE for each combination.

```

from sklearn.kernel_ridge import KernelRidge
from sklearn.metrics import mean_squared_error
import time
import pandas as pd # Import pandas to store results
from sklearn.pipeline import Pipeline # Import Pipeline

# Assuming X_train_reg, X_test_reg, y_train_reg, y_test_reg, and preprocessor_r

print("\nTraining Kernel Ridge Regression with Polynomial kernel and tuning hyp

# Define the hyperparameters to tune for the polynomial kernel
param_grid_poly = {
    'alpha': [0.1, 1.0, 10.0],
    'gamma': [0.01, 0.1, 1.0], # Using specific gamma values instead of 'auto'
    'degree': [2, 3] # Common degrees to test
}

results_kr_poly = []

# Perform manual grid search
for alpha_val in param_grid_poly['alpha']:
    for gamma_val in param_grid_poly['gamma']:
        for degree_val in param_grid_poly['degree']:
            print(f"\nTraining with alpha={alpha_val}, gamma={gamma_val}, degree={degree_val}")
            start_time_kr_poly = time.time()

            # Create a pipeline with preprocessing and Kernel Ridge (Polynomial)
            model_kr_poly = Pipeline(steps=[('preprocessor', preprocessor_reg),
                                             ('regressor', KernelRidge(kernel='poly'))])

            model_kr_poly.fit(X_train_reg, y_train_reg)
            end_time_kr_poly = time.time()
            training_time_kr_poly = end_time_kr_poly - start_time_kr_poly

            y_pred_kr_poly = model_kr_poly.predict(X_test_reg)
            mse_kr_poly = mean_squared_error(y_test_reg, y_pred_kr_poly)

            print(f" Training Time: {training_time_kr_poly:.4f} seconds")
            print(f" Validation MSE: {mse_kr_poly:.4f}")

            results_kr_poly.append({
                "alpha": alpha_val,
                "gamma": gamma_val,
                "degree": degree_val,
                "training_time": training_time_kr_poly,
                "validation_mse": mse_kr_poly
            })

# Store results in a DataFrame
df_kr_poly = pd.DataFrame(results_kr_poly)

```

```

print("\n--- Kernel Ridge (Polynomial) Results ---")
print(df_kr_poly.to_string(index=False))

# Find the best parameters and lowest MSE
best_kr_poly = df_kr_poly.loc[df_kr_poly['validation_mse'].idxmin()]

print(f"\nBest Kernel Ridge (Polynomial) Model:")
print(f"  Alpha: {best_kr_poly['alpha']} ")
print(f"  Gamma: {best_kr_poly['gamma']} ")
print(f"  Degree: {best_kr_poly['degree']} ")
print(f"  Lowest Validation MSE: {best_kr_poly['validation_mse']:.4f}")
print(f"  Training Time for Best Model: {best_kr_poly['training_time']:.4f} sec

# Add best result to the overall regression results dictionary
results_regression["Kernel Ridge (Poly)"] = {
    "Best Alpha": best_kr_poly['alpha'],
    "Best Gamma": best_kr_poly['gamma'],
    "Best Degree": best_kr_poly['degree'],
    "training_time": best_kr_poly['training_time'],
    "validation_mse": best_kr_poly['validation_mse']
}

```

Training Kernel Ridge Regression with Polynomial kernel and tuning hyperparameters

Training with alpha=0.1, gamma=0.01, degree=2

Training Time: 3.9927 seconds

Validation MSE: 2.9225

Training with alpha=0.1, gamma=0.01, degree=3

Training Time: 5.5008 seconds

Validation MSE: 2.9119

Training with alpha=0.1, gamma=0.1, degree=2

Training Time: 3.9656 seconds

Validation MSE: 3.2366

Training with alpha=0.1, gamma=0.1, degree=3

Training Time: 6.3505 seconds

Validation MSE: 3.8446

Training with alpha=0.1, gamma=1.0, degree=2

Training Time: 4.1135 seconds

Validation MSE: 4.9161

Training with alpha=0.1, gamma=1.0, degree=3

Training Time: 4.7253 seconds

Validation MSE: 14.2691

Training with alpha=1.0, gamma=0.01, degree=2

Training Time: 5.2973 seconds

Validation MSE: 3.0808

```
Training with alpha=1.0, gamma=0.01, degree=3
```

```
Training Time: 4.7207 seconds
```

```
Validation MSE: 3.0400
```

```
Training with alpha=1.0, gamma=0.1, degree=2
```

```
Training Time: 4.0406 seconds
```

```
Validation MSE: 2.9416
```

```
Training with alpha=1.0, gamma=0.1, degree=3
```

```
Training Time: 5.2194 seconds
```

```
Validation MSE: 3.0595
```

```
Training with alpha=1.0, gamma=1.0, degree=2
```

```
Training Time: 4.0210 seconds
```

```
Validation MSE: 4.2072
```

```
Training with alpha=1.0, gamma=1.0, degree=3
```

```
Training Time: 6.8057 seconds
```

```
Validation MSE: 7.3476
```

```
Training with alpha=10.0, gamma=0.01, degree=2
```

```
Training Time: 4.1043 seconds
```

```
Validation MSE: 3.2316
```

```
Training with alpha=10.0, gamma=0.01, degree=3
```

```
Training Time: 4.7817 seconds
```

```
Validation MSE: 3.2125
```

```
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression, Ridge
from sklearn.metrics import mean_squared_error
import time
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer # Import SimpleImputer

# Re-split the data specifically for regression features and target
# Use the X_reg and y_reg created in the previous step
X_train_reg, X_test_reg, y_train_reg, y_test_reg = train_test_split(
    X_reg, y_reg, test_size=0.2, random_state=42
)

print("Regression data split into training and testing sets.")
print(f"Training features shape: {X_train_reg.shape}")
print(f"Testing features shape: {X_test_reg.shape}")
print(f"Training target shape: {y_train_reg.shape}")
print(f"Testing target shape: {y_test_reg.shape}")

# Define the preprocessor again, ensuring it matches the features used for reg
# Use the categorical features reg and numerical features reg lists
```

```

preprocessor_reg = ColumnTransformer(
    transformers=[
        ('num', Pipeline([
            ('imputer', SimpleImputer(strategy='constant', fill_value=0)), # Im
            ('scaler', StandardScaler())
        ]), numerical_features_reg),
        ('cat', OneHotEncoder(handle_unknown='ignore')), categorical_features_re
    ],
    remainder='passthrough'
)

# --- Train and Evaluate OLS ---
print("\nTraining Ordinary Least Squares (OLS) model...")
start_time_lr = time.time()

# Create a pipeline that first preprocesses the data and then applies Linear Re
model_lr = Pipeline(steps=[('preprocessor', preprocessor_reg),
                           ('regressor', LinearRegression())])

model_lr.fit(X_train_reg, y_train_reg)
end_time_lr = time.time()
training_time_lr = end_time_lr - start_time_lr

y_pred_lr = model_lr.predict(X_test_reg)
mse_lr = mean_squared_error(y_test_reg, y_pred_lr)

print(f"OLS Training Time: {training_time_lr:.4f} seconds")
print(f"OLS Validation MSE: {mse_lr:.4f}")

# --- Train and Evaluate Ridge Regression ---
print("\nTraining Ridge Regression model...")
start_time_ridge = time.time()

# Create a pipeline for Ridge Regression
model_ridge = Pipeline(steps=[('preprocessor', preprocessor_reg),
                               ('regressor', Ridge(alpha=1.0))]) # Using default

model_ridge.fit(X_train_reg, y_train_reg)
end_time_ridge = time.time()
training_time_ridge = end_time_ridge - start_time_ridge

y_pred_ridge = model_ridge.predict(X_test_reg)
mse_ridge = mean_squared_error(y_test_reg, y_pred_ridge)

print(f'Ridge Training Time (alpha=1.0): {training_time_ridge:.4f} seconds')
print(f'Ridge Validation MSE (alpha=1.0): {mse_ridge:.4f}')

# Store results for later comparison

```

```

results_regression = {
    "OLS": {"training_time": training_time_lr, "validation_mse": mse_lr},
    "Ridge (alpha=1.0)": {"training_time": training_time_ridge, "validation_mse": mse_ridge}
}

Regression data split into training and testing sets.
Training features shape: (5281, 13)
Testing features shape: (1321, 13)
Training target shape: (5281,)
Testing target shape: (1321,)

Training Ordinary Least Squares (OLS) model...
OLS Training Time: 0.0781 seconds
OLS Validation MSE: 2.8933

Training Ridge Regression model...
Ridge Training Time (alpha=1.0): 0.0514 seconds
Ridge Validation MSE (alpha=1.0): 2.9055

```

Train and Evaluate Kernel Ridge Regression (Polynomial Kernel)

We'll continue exploring Kernel Ridge Regression, this time using a **Polynomial kernel**. This kernel allows the model to capture non-linear relationships by essentially creating new features that are combinations of the original features raised to different powers.

For the polynomial kernel, we'll tune three hyperparameters:

1. **alpha** : Similar to Ridge Regression, controls the regularization strength.
2. **gamma** : Influences the shape of the polynomial. If `gamma='auto'`, it will use $1/n_features$.
3. **degree** : The degree of the polynomial. A higher degree means more complex interactions between features.

We'll perform another manual grid search over these parameters and track the training time and validation MSE for each combination.

```

from sklearn.kernel_ridge import KernelRidge
from sklearn.metrics import mean_squared_error
import time
import pandas as pd # Import pandas to store results
from sklearn.pipeline import Pipeline # Import Pipeline

# Assuming X_train_reg, X_test_reg, y_train_reg, y_test_reg, and preprocessor_r

print("\nTraining Kernel Ridge Regression with Polynomial kernel and tuning hyp"

```

```

# Define the hyperparameters to tune for the polynomial kernel
param_grid_poly = {
    'alpha': [0.1, 1.0, 10.0],
    'gamma': [0.01, 0.1, 1.0], # Using specific gamma values instead of 'auto'
    'degree': [2, 3] # Common degrees to test
}

results_kr_poly = []

# Perform manual grid search
for alpha_val in param_grid_poly['alpha']:
    for gamma_val in param_grid_poly['gamma']:
        for degree_val in param_grid_poly['degree']:
            print(f"\nTraining with alpha={alpha_val}, gamma={gamma_val}, degree={degree_val}")
            start_time_kr_poly = time.time()

            # Create a pipeline with preprocessing and Kernel Ridge (Polynomial)
            model_kr_poly = Pipeline(steps=[('preprocessor', preprocessor_reg),
                                             ('regressor', KernelRidge(kernel='poly'))])

            model_kr_poly.fit(X_train_reg, y_train_reg)
            end_time_kr_poly = time.time()
            training_time_kr_poly = end_time_kr_poly - start_time_kr_poly

            y_pred_kr_poly = model_kr_poly.predict(X_test_reg)
            mse_kr_poly = mean_squared_error(y_test_reg, y_pred_kr_poly)

            print(f"  Training Time: {training_time_kr_poly:.4f} seconds")
            print(f"  Validation MSE: {mse_kr_poly:.4f}")

            results_kr_poly.append({
                "alpha": alpha_val,
                "gamma": gamma_val,
                "degree": degree_val,
                "training_time": training_time_kr_poly,
                "validation_mse": mse_kr_poly
            })

# Store results in a DataFrame
df_kr_poly = pd.DataFrame(results_kr_poly)

print("\n--- Kernel Ridge (Polynomial) Results ---")
print(df_kr_poly.to_string(index=False))

# Find the best parameters and lowest MSE
best_kr_poly = df_kr_poly.loc[df_kr_poly['validation_mse'].idxmin()]

print(f"\nBest Kernel Ridge (Polynomial) Model:")
print(f"  Alpha: {best_kr_poly['alpha']}"))
print(f"  Gamma: {best_kr_poly['gamma']}"))

```

```
print(f"  Degree: {best_kr_poly['degree']}")  
print(f"  Lowest Validation MSE: {best_kr_poly['validation_mse']:.4f}")  
print(f"  Training Time for Best Model: {best_kr_poly['training_time']:.4f} sec  
  
# Add best result to the overall regression results dictionary  
results_regression["Kernel Ridge (Poly)"] = {  
    "Best Alpha": best_kr_poly['alpha'],  
    "Best Gamma": best_kr_poly['gamma'],  
    "Best Degree": best_kr_poly['degree'],  
    "training_time": best_kr_poly['training_time'],  
    "validation_mse": best_kr_poly['validation_mse']  
}  
}
```

Training Kernel Ridge Regression with Polynomial kernel and tuning hyperparameters

Training with alpha=0.1, gamma=0.01, degree=2

Training Time: 4.3613 seconds

Validation MSE: 2.9532

Training with alpha=0.1, gamma=0.01, degree=3

Training Time: 5.9169 seconds

Validation MSE: 2.9821

Training with alpha=0.1, gamma=0.1, degree=2

Training Time: 4.4040 seconds

Validation MSE: 4.0337

Training with alpha=0.1, gamma=0.1, degree=3

Training Time: 6.9374 seconds

Validation MSE: 4.4556

Training with alpha=0.1, gamma=1.0, degree=2

Training Time: 4.4044 seconds

Validation MSE: 14.2084

Training with alpha=0.1, gamma=1.0, degree=3

Training Time: 5.2556 seconds

Validation MSE: 4.7054

Training with alpha=1.0, gamma=0.01, degree=2

Training Time: 4.9052 seconds

Validation MSE: 3.1066

Training with alpha=1.0, gamma=0.01, degree=3

Training Time: 5.0264 seconds

Validation MSE: 3.0627

Training with alpha=1.0, gamma=0.1, degree=2

Training Time: 5.7615 seconds

Validation MSE: 3.1341

Training with alpha=1.0, gamma=0.1, degree=3

Training Time: 5.0700 seconds

.....

```
Validation MSE: 3.5001

Training with alpha=1.0, gamma=1.0, degree=2
Training Time: 5.9905 seconds
Validation MSE: 6.5732

Training with alpha=1.0, gamma=1.0, degree=3
Training Time: 5.3198 seconds
Validation MSE: 4.6814

Training with alpha=10.0, gamma=0.01, degree=2
Training Time: 4.3967 seconds
Validation MSE: 3.2721

Training with alpha=10.0, gamma=0.01, degree=3
Training Time: 5.7476 seconds
Validation MSE: 3.2439
```

▼ Train and Evaluate OLS and Ridge Regression

We're starting with two basic but important types of regression:

1. **Ordinary Least Squares (OLS)**: This is like trying to find the straight line that best fits our data points to predict the goal difference. It's a fundamental method that's easy to understand.
 2. **Ridge Regression**: This is similar to OLS but with a slight twist. It adds a penalty to the model to prevent it from becoming too complex or overly sensitive to individual data points. This can sometimes help the model perform better on new, unseen data.

We'll train both models and see how well they predict the goal difference on our testing data. We'll measure their performance using something called **Mean Squared Error (MSE)**. A lower MSE means the model's predictions are closer to the actual goal differences. We'll also keep track of how long it takes to train each model.

```
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression, Ridge
from sklearn.metrics import mean_squared_error
import time
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer # Import SimpleImputer

# Re-split the data specifically for regression features and target
# Use the X_reg and y_reg created in the previous step
```

```

x_train_reg, x_test_reg, y_train_reg, y_test_reg = train_test_split(
    X_reg, y_reg, test_size=0.2, random_state=42
)

print("Regression data split into training and testing sets.")
print(f"Training features shape: {X_train_reg.shape}")
print(f"Testing features shape: {X_test_reg.shape}")
print(f"Training target shape: {y_train_reg.shape}")
print(f"Testing target shape: {y_test_reg.shape}")

# Define the preprocessor again, ensuring it matches the features used for regr
# Use the categorical_features_reg and numerical_features_reg lists
preprocessor_reg = ColumnTransformer(
    transformers=[
        ('num', Pipeline([
            ('imputer', SimpleImputer(strategy='constant', fill_value=0)), # Im
            ('scaler', StandardScaler())
        ]), numerical_features_reg),
        ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features_re
    ],
    remainder='passthrough'
)

# --- Train and Evaluate OLS ---
print("\nTraining Ordinary Least Squares (OLS) model...")
start_time_lr = time.time()

# Create a pipeline that first preprocesses the data and then applies Linear Re
model_lr = Pipeline(steps=[('preprocessor', preprocessor_reg),
                           ('regressor', LinearRegression())])

model_lr.fit(X_train_reg, y_train_reg)
end_time_lr = time.time()
training_time_lr = end_time_lr - start_time_lr

y_pred_lr = model_lr.predict(X_test_reg)
mse_lr = mean_squared_error(y_test_reg, y_pred_lr)

print(f"OLS Training Time: {training_time_lr:.4f} seconds")
print(f"OLS Validation MSE: {mse_lr:.4f}")

# --- Train and Evaluate Ridge Regression ---
print("\nTraining Ridge Regression model...")
start_time_ridge = time.time()

# Create a pipeline for Ridge Regression
model_ridge = Pipeline(steps=[('preprocessor', preprocessor_reg),
                             ('regressor', Ridge(alpha=1.0))]) # Using default

```

```

model_ridge.fit(X_train_reg, y_train_reg)
end_time_ridge = time.time()
training_time_ridge = end_time_ridge - start_time_ridge

y_pred_ridge = model_ridge.predict(X_test_reg)
mse_ridge = mean_squared_error(y_test_reg, y_pred_ridge)

print(f"Ridge Training Time (alpha=1.0): {training_time_ridge:.4f} seconds")
print(f"Ridge Validation MSE (alpha=1.0): {mse_ridge:.4f}")

# Store results for later comparison
results_regression = {
    "OLS": {"training_time": training_time_lr, "validation_mse": mse_lr},
    "Ridge (alpha=1.0)": {"training_time": training_time_ridge, "validation_mse": mse_ridge}
}

```

Regression data split into training and testing sets.

Training features shape: (5281, 13)

Testing features shape: (1321, 13)

Training target shape: (5281,)

Testing target shape: (1321,)

Training Ordinary Least Squares (OLS) model...

OLS Training Time: 0.0714 seconds

OLS Validation MSE: 2.8933

Training Ridge Regression model...

Ridge Training Time (alpha=1.0): 0.0518 seconds

Ridge Validation MSE (alpha=1.0): 2.9055

```

# Assuming 'cleaned_df' is the DataFrame with engineered features from previous

# Define the features (X) - these are the columns we'll use to predict the goal
# We'll exclude the actual goal difference and individual scores from the features
# as those would leak information about the outcome. We'll include the engineered
# features like 'expected_goal_diff', date/time features, and team/competition identifiers

# Let's select the relevant features that would be available *before* the match
# We'll use the categorical features and the engineered numerical features like
categorical_features_reg = [
    'competition_name', 'competition_code',
    'homeTeam_name', 'homeTeam_tla',
    'awayTeam_name', 'awayTeam_tla'
]

numerical_features_reg = [
    'expected_goal_diff', # This is our pre-match proxy for goal difference
    'match_hour',
    'match_day_of_week',
    'match_month',

```

```
'is_home_game', # This is our binary home advantage indicator
'homeTeam_rolling_avg_goals_scored_5', # Rolling averages from previous cod
'awayTeam_rolling_avg_goals_conceded_5' # Rolling averages from previous cc
]

# Filter out any features that might not exist in the DataFrame based on previous
numerical_features_reg = [f for f in numerical_features_reg if f in cleaned_df]
categorical_features_reg = [f for f in categorical_features_reg if f in cleaned_df]

# Combine the feature lists
all_features_reg = categorical_features_reg + numerical_features_reg

# Create the features DataFrame
X_reg = cleaned_df[all_features_reg].copy()

# Define the target variable (y) - this is what we want to predict.
y_reg = cleaned_df['goal_difference']

print("Features (X_reg) and target (y_reg) for regression defined.")
print(f"Features shape: {X_reg.shape}")
print(f"Target shape: {y_reg.shape}")
print("\nFirst 5 rows of features:")
display(X_reg.head())
print("\nFirst 5 rows of target:")
display(y_reg.head())
```

```
import requests
import pandas as pd
import time

# YOUR API KEY GOES HERE
api_key = "e0d3c18b3e2c496e9a79bf92ccbbe53f"
BASE_URL = "https://api.football-data.org/v4/"

headers = {
    'X-Auth-Token': api_key,
    'X-Unfold-Goals': 'true'
}

# The leagues to include, based on your free plan
competition_codes = ['BL1', 'DED', 'BSA', 'PD', 'FL1', 'ELC', 'PPL', 'SA', 'PL']
seasons = [2023, 2024]

all_matches = []
unique_match_ids = set()

print("Starting data collection for multiple leagues and seasons...")
for competition_code in competition_codes:
    for season in seasons:
        print(f"\nFetching matches for {competition_code} in {season}/{season+1}")
        url = f"{BASE_URL}competitions/{competition_code}/matches?season={seasc

try:
    response = requests.get(url, headers=headers)
    response.raise_for_status()
    data = response.json()

    if 'matches' in data:
        for match in data['matches']:
            if match['id'] not in unique_match_ids:
                all_matches.append(match)
                unique_match_ids.add(match['id'])
                print(f" Successfully fetched {len(data['matches'])} matches.")
            else:
                print(f" Match ID {match['id']} already exists in the dataset.")

    # Save the data to a CSV file
    df = pd.DataFrame(all_matches)
    df.to_csv('football_data.csv', index=False)
    print("Data collection completed and saved to football_data.csv")
```

```
        else:
            print(f"  No matches found for {competition_code} in that seasc

        except requests.exceptions.RequestException as e:
            print(f"  An error occurred: {e}")

        time.sleep(15) # Respect the API rate limit

    # --- Create a Final DataFrame ---
    print("\nCreating a single DataFrame from the collected data...")
    if all_matches:
        final_df = pd.DataFrame(all_matches)

        print(f"\nTotal unique matches in the final DataFrame: {len(final_df)}")

        print("\nDataFrame created. Here are the first 5 rows:")
        print(final_df.head())

        print("\nDataFrame Columns:")
        print(final_df.columns)
    else:
        print("No match data was collected.")
```

Starting data collection for multiple leagues and seasons...

Fetching matches for BL1 in 2023/2024 season...
Successfully fetched 306 matches.

Fetching matches for BL1 in 2024/2025 season...
Successfully fetched 306 matches.

Fetching matches for DED in 2023/2024 season...
Successfully fetched 306 matches.

Fetching matches for DED in 2024/2025 season...
Successfully fetched 306 matches.

Fetching matches for BSA in 2023/2024 season...
Successfully fetched 380 matches.

Fetching matches for BSA in 2024/2025 season...
Successfully fetched 380 matches.

Fetching matches for PD in 2023/2024 season...
Successfully fetched 380 matches.

Fetching matches for PD in 2024/2025 season...
Successfully fetched 380 matches.

Fetching matches for FL1 in 2023/2024 season...
Successfully fetched 306 matches.

Fetching matches for FL1 in 2024/2025 season...
Successfully fetched 306 matches

```
Successfully fetched 500 matches.
```

```
Fetching matches for ELC in 2023/2024 season...
Successfully fetched 557 matches.
```

```
Fetching matches for ELC in 2024/2025 season...
Successfully fetched 557 matches.
```

```
Fetching matches for PPL in 2023/2024 season...
Successfully fetched 306 matches.
```

```
Fetching matches for PPL in 2024/2025 season...
Successfully fetched 306 matches.
```

```
Fetching matches for SA in 2023/2024 season...
Successfully fetched 380 matches.
```

```
Fetching matches for SA in 2024/2025 season...
Successfully fetched 380 matches.
```

```
Fetching matches for PL in 2023/2024 season...
Successfully fetched 380 matches.
```

```
Fetching matches for PL in 2024/2025 season...
Successfully fetched 380 matches.
```

```
Creating a single DataFrame from the collected data...
```

```
# --- Step 1: Flattening the Data ---
# Extract key information from the nested columns.

# Normalize home and away team data
home_teams = pd.json_normalize(final_df['homeTeam']).add_prefix('homeTeam_')
away_teams = pd.json_normalize(final_df['awayTeam']).add_prefix('awayTeam_')

# Normalize score data
scores = pd.json_normalize(final_df['score'])

# Normalize competition data to get the league name
competitions = pd.json_normalize(final_df['competition']).add_prefix('competiti')

# Combine all the flattened data into a new DataFrame
cleaned_df = pd.concat([
    final_df[['id', 'utcDate', 'status']], # Keep essential columns
    competitions[['competition_name', 'competition_code']],
    home_teams[['homeTeam_name', 'homeTeam_tla']],
    away_teams[['awayTeam_name', 'awayTeam_tla']],
    scores[['winner', 'fullTime.home', 'fullTime.away', 'halfTime.home', 'halfT'],
], axis=1)

# --- Step 2: Handle Missing Values ---
# Fill the few missing halftime scores with 0
```

```

cleaned_df['halfTime.home'] = cleaned_df['halfTime.home'].fillna(0).astype(int)
cleaned_df['halfTime.away'] = cleaned_df['halfTime.away'].fillna(0).astype(int)

# Drop columns with too many missing values
columns_to_drop = ['regularTime.home', 'regularTime.away', 'extraTime.home', 'e
cleaned_df.drop(columns=columns_to_drop, inplace=True, errors='ignore')

# --- Step 3: Rename Columns for Clarity ---
cleaned_df.rename(columns={
    'fullTime.home': 'home_score',
    'fullTime.away': 'away_score',
    'halfTime.home': 'halfTime_home_score',
    'halfTime.away': 'halfTime_away_score'
}, inplace=True)

print("Comprehensive cleaning complete.\n")
print(f"DataFrame now has {len(cleaned_df.columns)} columns and {len(cleaned_df)} rows")
print("The first 5 rows of the new, clean DataFrame:")
print(cleaned_df.head())
print("\nNew DataFrame Columns:")
print(cleaned_df.columns)

```

Comprehensive cleaning complete.

DataFrame now has 14 columns and 6602 rows.

The first 5 rows of the new, clean DataFrame:

	id	utcDate	status	competition_name	competition_code	
0	441789	2023-08-18T18:30:00Z	FINISHED	Bundesliga	BL1	
1	441792	2023-08-19T13:30:00Z	FINISHED	Bundesliga	BL1	
2	441794	2023-08-19T13:30:00Z	FINISHED	Bundesliga	BL1	
3	441795	2023-08-19T13:30:00Z	FINISHED	Bundesliga	BL1	
4	441796	2023-08-19T13:30:00Z	FINISHED	Bundesliga	BL1	

	homeTeam_name	homeTeam_tla	awayTeam_name	awayTeam_tla	
0	SV Werder Bremen	SVW	FC Bayern München	FCB	
1	Bayer 04 Leverkusen	B04	RB Leipzig	RBL	
2	VfL Wolfsburg	WOB	1. FC Heidenheim 1846	HEI	
3	TSG 1899 Hoffenheim	TSG	SC Freiburg	SCF	
4	FC Augsburg	FCA	Borussia Mönchengladbach	BMG	

	winner	home_score	away_score	halfTime_home_score	halfTime_away_score	
0	AWAY_TEAM	0	4	0	1	
1	HOME_TEAM	3	2	2	1	
2	HOME_TEAM	2	0	2	0	
3	AWAY_TEAM	1	2	0	2	
4	DRAW	4	4	3	3	

New DataFrame Columns:

```

Index(['id', 'utcDate', 'status', 'competition_name', 'competition_code',
       'homeTeam_name', 'homeTeam_tla', 'awayTeam_name', 'awayTeam_tla',
       'winner', 'home_score', 'away_score', 'halfTime_home_score',
       'halfTime_away_score'],
      dtype='object')

```

```

# --- Create safe pre-match proxy for goal_difference ---
# Calculate per-match goal differences (post-match, but we'll use them only for
cleaned_df['home_goal_diff'] = cleaned_df['home_score'] - cleaned_df['away_score']
cleaned_df['away_goal_diff'] = cleaned_df['away_score'] - cleaned_df['home_score']

# Sort by date so rolling windows are correct
cleaned_df = cleaned_df.sort_values(by='utcDate')

# Rolling averages for last 5 matches (excluding current match)
cleaned_df['home_avg_goal_diff_last5'] = (
    cleaned_df.groupby('homeTeam_name')['home_goal_diff']
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)

cleaned_df['away_avg_goal_diff_last5'] = (
    cleaned_df.groupby('awayTeam_name')['away_goal_diff']
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)

# Combine into a single "expected goal diff" feature
cleaned_df['expected_goal_diff'] = (
    cleaned_df['home_avg_goal_diff_last5'] - cleaned_df['away_avg_goal_diff_last5']
)

# Convert the 'utcDate' column to datetime objects
cleaned_df['utcDate'] = pd.to_datetime(cleaned_df['utcDate'])

# Extract date and time components
cleaned_df['match_hour'] = cleaned_df['utcDate'].dt.hour
cleaned_df['match_day_of_week'] = cleaned_df['utcDate'].dt.dayofweek # Monday=0
cleaned_df['match_month'] = cleaned_df['utcDate'].dt.month

# Create a binary feature for home advantage (1 if home game, 0 otherwise)
cleaned_df['is_home_game'] = 1 # In this dataset structure, each row is a match

# Add rolling averages for goals scored and conceded for home and away teams
# (Assuming these were intended to be created from previous steps based on the
cleaned_df['homeTeam_rolling_avg_goals_scored_5'] = (
    cleaned_df.groupby('homeTeam_name')['home_score']
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)

cleaned_df['awayTeam_rolling_avg_goals_conceded_5'] = (
    cleaned_df.groupby('awayTeam_name')['home_score'] # Away team conceded is home score
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)

print("\nEngineered features added to cleaned_df.\n")

```

```
display(cleaned_df.head())
```

```
# Define the features (X) - these are the columns we'll use to predict the goal
# We'll exclude the actual goal difference and individual scores from the features
# as those would leak information about the outcome. We'll include the engineered
# features like 'expected_goal_diff', date/time features, and team/competition identifiers
# from the cleaned DataFrame.
cleaned_df['goal_difference'] = cleaned_df['home_score'] - cleaned_df['away_score']

# Let's select the relevant features that would be available *before* the match
# We'll use the categorical features and the engineered numerical features like
categorical_features_reg = [
    'competition_name', 'competition_code',
    'homeTeam_name', 'homeTeam_tla',
    'awayTeam_name', 'awayTeam_tla'
]

numerical_features_reg = [
    'expected_goal_diff', # This is our pre-match proxy for goal difference
    'match_hour',
    'match_day_of_week',
    'match_month',
    'is_home_game', # This is our binary home advantage indicator
    'homeTeam_rolling_avg_goals_scored_5', # Rolling averages from previous 5 games
    'awayTeam_rolling_avg_goals_conceded_5' # Rolling averages from previous 5 games
]

# Filter out any features that might not exist in the DataFrame based on previous definitions
numerical_features_reg = [f for f in numerical_features_reg if f in cleaned_df]
categorical_features_reg = [f for f in categorical_features_reg if f in cleaned_df]
```

```
# Combine the feature lists
all_features_reg = categorical_features_reg + numerical_features_reg

# Create the features DataFrame
X_reg = cleaned_df[all_features_reg].copy()

# Define the target variable (y) - this is what we want to predict.
y_reg = cleaned_df['goal_difference']

print("Features (X_reg) and target (y_reg) for regression defined.")
print(f"Features shape: {X_reg.shape}")
print(f"Target shape: {y_reg.shape}")
print("\nFirst 5 rows of features:")
display(X_reg.head())
print("\nFirst 5 rows of target:")
display(y_reg.head())
```

✓ Prepare Data for Regression

Before we train our regression models, we need to make sure our data is in the right shape and format. Think of this like getting your ingredients ready before you start cooking!

For regression, which is about predicting a continuous number (like the difference in scores), our "target" (what we want to predict) needs to be a single number for each match. We already created a 'goal_difference' column earlier, which is perfect for this. It's the home team's score minus the away team's score.

Our "features" (the information we use to make the prediction) are all the other columns we've cleaned and engineered. We used scaling on these features before for classification, which helps models work better, and that scaling is still appropriate for most regression models.

So, in this step, we'll just explicitly define our features (X) and our numerical target (y) for the regression task.

```
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression, Ridge
from sklearn.metrics import mean_squared_error
import time
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer # Import SimpleImputer

# Re-split the data specifically for regression features and target
# Use the X_reg and y_reg created in the previous step
X_train_reg, X_test_reg, y_train_reg, y_test_reg = train_test_split(
    X_reg, y_reg, test_size=0.2, random_state=42
)

print("Regression data split into training and testing sets.")
print(f"Training features shape: {X_train_reg.shape}")
print(f"Testing features shape: {X_test_reg.shape}")
print(f"Training target shape: {y_train_reg.shape}")
print(f"Testing target shape: {y_test_reg.shape}")
```

```

# Define the preprocessor again, ensuring it matches the features used for regr
# Use the categorical_features_reg and numerical_features_reg lists
preprocessor_reg = ColumnTransformer(
    transformers=[
        ('num', Pipeline([
            ('imputer', SimpleImputer(strategy='constant', fill_value=0)), # Im
            ('scaler', StandardScaler())
        ]), numerical_features_reg),
        ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features_re
    ],
    remainder='passthrough'
)

# --- Train and Evaluate OLS ---
print("\nTraining Ordinary Least Squares (OLS) model...")
start_time_lr = time.time()

# Create a pipeline that first preprocesses the data and then applies Linear Re
model_lr = Pipeline(steps=[('preprocessor', preprocessor_reg),
                           ('regressor', LinearRegression())])

model_lr.fit(X_train_reg, y_train_reg)
end_time_lr = time.time()
training_time_lr = end_time_lr - start_time_lr

y_pred_lr = model_lr.predict(X_test_reg)
mse_lr = mean_squared_error(y_test_reg, y_pred_lr)

print(f"OLS Training Time: {training_time_lr:.4f} seconds")
print(f"OLS Validation MSE: {mse_lr:.4f}")

# --- Train and Evaluate Ridge Regression ---
print("\nTraining Ridge Regression model...")
start_time_ridge = time.time()

# Create a pipeline for Ridge Regression
model_ridge = Pipeline(steps=[('preprocessor', preprocessor_reg),
                             ('regressor', Ridge(alpha=1.0))]) # Using default

model_ridge.fit(X_train_reg, y_train_reg)
end_time_ridge = time.time()
training_time_ridge = end_time_ridge - start_time_ridge

y_pred_ridge = model_ridge.predict(X_test_reg)
mse_ridge = mean_squared_error(y_test_reg, y_pred_ridge)

print(f'Ridge Training Time (alpha=1.0): {training_time_ridge:.4f} seconds')
print(f'Ridge Validation MSE (alpha=1.0): {mse_ridge:.4f}')

```

```

# Store results for later comparison
results_regression = {
    "OLS": {"training_time": training_time_lr, "validation_mse": mse_lr},
    "Ridge (alpha=1.0)": {"training_time": training_time_ridge, "validation_mse": mse_ridge}
}

Regression data split into training and testing sets.
Training features shape: (5281, 13)
Testing features shape: (1321, 13)
Training target shape: (5281,)
Testing target shape: (1321,)

Training Ordinary Least Squares (OLS) model...
OLS Training Time: 0.0633 seconds
OLS Validation MSE: 2.8782

Training Ridge Regression model...
Ridge Training Time (alpha=1.0): 0.0634 seconds
Ridge Validation MSE (alpha=1.0): 2.8870

```

```

# Assuming 'cleaned_df' is the DataFrame with engineered features from previous

# Define the features (X) - these are the columns we'll use to predict the goal
# We'll exclude the actual goal difference and individual scores from the features
# as those would leak information about the outcome. We'll include the engineered
# like 'expected_goal_diff', date/time features, and team/competition identifiers

# Let's select the relevant features that would be available *before* the match
# We'll use the categorical features and the engineered numerical features like
categorical_features_reg = [
    'competition_name', 'competition_code',
    'homeTeam_name', 'homeTeam_tla',
    'awayTeam_name', 'awayTeam_tla'
]

numerical_features_reg = [
    'expected_goal_diff', # This is our pre-match proxy for goal difference
    'match_hour',
    'match_day_of_week',
    'match_month',
    'is_home_game', # This is our binary home advantage indicator
    'homeTeam_rolling_avg_goals_scored_5', # Rolling averages from previous code
    'awayTeam_rolling_avg_goals_conceded_5' # Rolling averages from previous code
]

# Filter out any features that might not exist in the DataFrame based on previous
numerical_features_reg = [f for f in numerical_features_reg if f in cleaned_df]
categorical_features_reg = [f for f in categorical_features_reg if f in cleaned_df]

```

```
# Combine the feature lists
all_features_reg = categorical_features_reg + numerical_features_reg

# Create the features DataFrame
X_reg = cleaned_df[all_features_reg].copy()

# Define the target variable (y) - this is what we want to predict.
y_reg = cleaned_df['goal_difference']

print("Features (X_reg) and target (y_reg) for regression defined.")
print(f"Features shape: {X_reg.shape}")
print(f"Target shape: {y_reg.shape}")
print("\nFirst 5 rows of features:")
display(X_reg.head())
print("\nFirst 5 rows of target:")
display(y_reg.head())
```

✓ Prepare Data for Regression

Before we train our regression models, we need to make sure our data is in the right shape and format. Think of this like getting your ingredients ready before you start cooking!

For regression, which is about predicting a continuous number (like the difference in scores), our "target" (what we want to predict) needs to be a single number for each match. We already created a 'goal_difference' column earlier, which is perfect for this. It's the home team's score minus the away team's score.

Our "features" (the information we use to make the prediction) are all the other columns we've cleaned and engineered. We used scaling on these features before for classification, which helps models work better, and that scaling is still appropriate for most regression models.

So, in this step, we'll just explicitly define our features (X) and our numerical target (y) for the regression task.

```
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression, Ridge
from sklearn.metrics import mean_squared_error
import time
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer # Import SimpleImputer

# Re-split the data specifically for regression features and target
# Use the X_reg and y_reg created in the previous step
X_train_reg, X_test_reg, y_train_reg, y_test_reg = train_test_split(
    X_reg, y_reg, test_size=0.2, random_state=42
)

print("Regression data split into training and testing sets.")
print(f"Training features shape: {X_train_reg.shape}")
print(f"Testing features shape: {X_test_reg.shape}")
print(f"Training target shape: {y_train_reg.shape}")
print(f"Testing target shape: {y_test_reg.shape}")

# Define the preprocessor again, ensuring it matches the features used for regr
```

```

# Use the categorical_features_reg and numerical_features_reg lists
preprocessor_reg = ColumnTransformer(
    transformers=[
        ('num', Pipeline([
            ('imputer', SimpleImputer(strategy='constant', fill_value=0)), # Im
            ('scaler', StandardScaler())
        ]), numerical_features_reg),
        ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features_re
    ],
    remainder='passthrough'
)

# --- Train and Evaluate OLS ---
print("\nTraining Ordinary Least Squares (OLS) model...")
start_time_lr = time.time()

# Create a pipeline that first preprocesses the data and then applies Linear Re
model_lr = Pipeline(steps=[('preprocessor', preprocessor_reg),
                           ('regressor', LinearRegression())])

model_lr.fit(X_train_reg, y_train_reg)
end_time_lr = time.time()
training_time_lr = end_time_lr - start_time_lr

y_pred_lr = model_lr.predict(X_test_reg)
mse_lr = mean_squared_error(y_test_reg, y_pred_lr)

print(f"OLS Training Time: {training_time_lr:.4f} seconds")
print(f"OLS Validation MSE: {mse_lr:.4f}")

# --- Train and Evaluate Ridge Regression ---
print("\nTraining Ridge Regression model...")
start_time_ridge = time.time()

# Create a pipeline for Ridge Regression
model_ridge = Pipeline(steps=[('preprocessor', preprocessor_reg),
                               ('regressor', Ridge(alpha=1.0))]) # Using default

model_ridge.fit(X_train_reg, y_train_reg)
end_time_ridge = time.time()
training_time_ridge = end_time_ridge - start_time_ridge

y_pred_ridge = model_ridge.predict(X_test_reg)
mse_ridge = mean_squared_error(y_test_reg, y_pred_ridge)

print(f'Ridge Training Time (alpha=1.0): {training_time_ridge:.4f} seconds')
print(f'Ridge Validation MSE (alpha=1.0): {mse_ridge:.4f}')

# Extra results for later comparison

```

```

# Score results for later comparison
results_regression = {
    "OLS": {"training_time": training_time_lr, "validation_mse": mse_lr},
    "Ridge (alpha=1.0)": {"training_time": training_time_ridge, "validation_mse": mse_ridge}
}

Regression data split into training and testing sets.
Training features shape: (5281, 13)
Testing features shape: (1321, 13)
Training target shape: (5281,)
Testing target shape: (1321,)

Training Ordinary Least Squares (OLS) model...
OLS Training Time: 0.0589 seconds
OLS Validation MSE: 2.8782

Training Ridge Regression model...
Ridge Training Time (alpha=1.0): 0.0583 seconds
Ridge Validation MSE (alpha=1.0): 2.8870

```

```

# --- Create safe pre-match proxy for goal_difference ---
# Calculate per-match goal differences (post-match, but we'll use them only for
# cleaned_df['home_goal_diff'] = cleaned_df['home_score'] - cleaned_df['away_score']
cleaned_df['away_goal_diff'] = cleaned_df['away_score'] - cleaned_df['home_score']

# Calculate the actual goal difference (home - away) for the target variable
cleaned_df['goal_difference'] = cleaned_df['home_score'] - cleaned_df['away_score']

# Sort by date so rolling windows are correct
cleaned_df = cleaned_df.sort_values(by='utcDate')

# Rolling averages for last 5 matches (excluding current match)
cleaned_df['home_avg_goal_diff_last5'] = (
    cleaned_df.groupby('homeTeam_name')['home_goal_diff']
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)

cleaned_df['away_avg_goal_diff_last5'] = (
    cleaned_df.groupby('awayTeam_name')['away_goal_diff']
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)

# Combine into a single "expected goal diff" feature
cleaned_df['expected_goal_diff'] = (
    cleaned_df['home_avg_goal_diff_last5'] - cleaned_df['away_avg_goal_diff_last5']
)

# Convert the 'utcDate' column to datetime objects
cleaned_df['utcDate'] = pd.to_datetime(cleaned_df['utcDate'])

# Extract date and time components

```

```

# EXTRACT DATE AND TIME COMPONENTS
cleaned_df['match_hour'] = cleaned_df['utcDate'].dt.hour
cleaned_df['match_day_of_week'] = cleaned_df['utcDate'].dt.dayofweek # Monday=0
cleaned_df['match_month'] = cleaned_df['utcDate'].dt.month

# Create a binary feature for home advantage (1 if home game, 0 otherwise)
cleaned_df['is_home_game'] = 1 # In this dataset structure, each row is a match

# Add rolling averages for goals scored and conceded for home and away teams
# (Assuming these were intended to be created from previous steps based on the
cleaned_df['homeTeam_rolling_avg_goals_scored_5'] = (
    cleaned_df.groupby('homeTeam_name')['home_score']
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)

cleaned_df['awayTeam_rolling_avg_goals_conceded_5'] = (
    cleaned_df.groupby('awayTeam_name')['home_score'] # Away team conceded is home score
    .transform(lambda x: x.shift().rolling(5, min_periods=1).mean())
)

print("\nEngineered features added to cleaned_df.")
display(cleaned_df.head())

```

```

# Define the features (X) - these are the columns we'll use to predict the goal
# We'll exclude the actual goal difference and individual scores from the features
# as those would leak information about the outcome. We'll include the engineered
# features like 'expected_goal_diff', date/time features, and team/competition identifiers

# Let's select the relevant features that would be available *before* the match
# We'll use the categorical features and the engineered numerical features like

```

```
# we'll use the categorical features and the engineered numerical features like
categorical_features_reg = [
    'competition_name', 'competition_code',
    'homeTeam_name', 'homeTeam_tla',
    'awayTeam_name', 'awayTeam_tla'
]

numerical_features_reg = [
    'expected_goal_diff', # This is our pre-match proxy for goal difference
    'match_hour',
    'match_day_of_week',
    'match_month',
    'is_home_game', # This is our binary home advantage indicator
    'homeTeam_rolling_avg_goals_scored_5', # Rolling averages from previous cod
    'awayTeam_rolling_avg_goals_conceded_5' # Rolling averages from previous cc
]

# Filter out any features that might not exist in the DataFrame based on previous
numerical_features_reg = [f for f in numerical_features_reg if f in cleaned_df]
categorical_features_reg = [f for f in categorical_features_reg if f in cleaned_df]

# Combine the feature lists
all_features_reg = categorical_features_reg + numerical_features_reg

# Create the features DataFrame
X_reg = cleaned_df[all_features_reg].copy()

# Define the target variable (y) - this is what we want to predict.
y_reg = cleaned_df['goal_difference']

print("Features (X_reg) and target (y_reg) for regression defined.")
print(f"Features shape: {X_reg.shape}")
print(f"Target shape: {y_reg.shape}")
print("\nFirst 5 rows of features:")
display(X_reg.head())
print("\nFirst 5 rows of target:")
display(y_reg.head())
```

✓ Prepare Data for Regression

Before we train our regression models, we need to make sure our data is in the right shape and format. Think of this like getting your ingredients ready before you start cooking!

For regression, which is about predicting a continuous number (like the difference in scores), our "target" (what we want to predict) needs to be a single number for each match. We already created a 'goal_difference' column earlier, which is perfect for this. It's the home team's score minus the away team's score.

Our "features" (the information we use to make the prediction) are all the other columns we've cleaned and engineered. We used scaling on these features before for classification, which helps models work better, and that scaling is still appropriate for most regression models.

So, in this step, we'll just explicitly define our features (X) and our numerical target (y) for the regression task.

```
from sklearn.model_selection import train_test_split  
from sklearn.linear_model import LinearRegression, Ridge  
from sklearn.metrics import mean_squared_error
```

```

from sklearn.metrics import mean_squared_error
import time
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer # Import SimpleImputer

# Re-split the data specifically for regression features and target
# Use the X_reg and y_reg created in the previous step
X_train_reg, X_test_reg, y_train_reg, y_test_reg = train_test_split(
    X_reg, y_reg, test_size=0.2, random_state=42
)

print("Regression data split into training and testing sets.")
print(f"Training features shape: {X_train_reg.shape}")
print(f"Testing features shape: {X_test_reg.shape}")
print(f"Training target shape: {y_train_reg.shape}")
print(f"Testing target shape: {y_test_reg.shape}")

# Define the preprocessor again, ensuring it matches the features used for reg
# Use the categorical_features_reg and numerical_features_reg lists
preprocessor_reg = ColumnTransformer(
    transformers=[
        ('num', Pipeline([
            ('imputer', SimpleImputer(strategy='constant', fill_value=0)), # Im
            ('scaler', StandardScaler())
        ]), numerical_features_reg),
        ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features_re
    ],
    remainder='passthrough'
)

# --- Train and Evaluate OLS ---
print("\nTraining Ordinary Least Squares (OLS) model...")
start_time_lr = time.time()

# Create a pipeline that first preprocesses the data and then applies Linear Re
model_lr = Pipeline(steps=[('preprocessor', preprocessor_reg),
                           ('regressor', LinearRegression())])

model_lr.fit(X_train_reg, y_train_reg)
end_time_lr = time.time()
training_time_lr = end_time_lr - start_time_lr

y_pred_lr = model_lr.predict(X_test_reg)
mse_lr = mean_squared_error(y_test_reg, y_pred_lr)

print(f"OLS Training Time: {training_time_lr:.4f} seconds")

```

```

print(f"OLS Validation MSE: {mse_lr:.4f}")

# --- Train and Evaluate Ridge Regression ---
print("\nTraining Ridge Regression model...")
start_time_ridge = time.time()

# Create a pipeline for Ridge Regression
model_ridge = Pipeline(steps=[('preprocessor', preprocessor_reg),
                             ('regressor', Ridge(alpha=1.0))]) # Using default

model_ridge.fit(X_train_reg, y_train_reg)
end_time_ridge = time.time()
training_time_ridge = end_time_ridge - start_time_ridge

y_pred_ridge = model_ridge.predict(X_test_reg)
mse_ridge = mean_squared_error(y_test_reg, y_pred_ridge)

print(f"Ridge Training Time (alpha=1.0): {training_time_ridge:.4f} seconds")
print(f"Ridge Validation MSE (alpha=1.0): {mse_ridge:.4f}")

# Store results for later comparison
results_regression = {
    "OLS": {"training_time": training_time_lr, "validation_mse": mse_lr},
    "Ridge (alpha=1.0)": {"training_time": training_time_ridge, "validation_mse": mse_ridge}
}

```

Regression data split into training and testing sets.

Training features shape: (5281, 13)

Testing features shape: (1321, 13)

Training target shape: (5281,)

Testing target shape: (1321,)

Training Ordinary Least Squares (OLS) model...

OLS Training Time: 0.0814 seconds

OLS Validation MSE: 2.8272

Training Ridge Regression model...

Ridge Training Time (alpha=1.0): 0.0902 seconds

Ridge Validation MSE (alpha=1.0): 2.8388

Train and Evaluate Shallow Neural Network (Baseline) for Soccer Data

Now that we have defined the architecture for the shallow baseline model, we need to train it on our processed soccer training data and evaluate its performance on the testing data.

```

import torch
import torch.nn as nn
import torch.optim as optim
from torch.utils.data import DataLoader # Import DataLoader if not already imported
import time # Import time

# Assuming shallow_baseline_model is defined from the previous step
# Assuming train_dataloader and test_dataloader are defined from the data preparation step

print("--- Starting Training for Shallow Baseline Neural Network ---")

# Define Loss Function and Optimizer
# For multi-class classification with one-hot encoded labels, CrossEntropyLoss is appropriate
# It combines LogSoftmax and NLLLoss. Note: our model's output layer does NOT have a softmax
# which is correct when using CrossEntropyLoss.
criterion = nn.CrossEntropyLoss()
# Adam optimizer is a good default choice
optimizer = optim.Adam(shallow_baseline_model.parameters(), lr=0.001) # Learning rate

# --- Training Loop ---
num_epochs = 50 # Define the number of training epochs (can be adjusted)
train_loss_history = []
val_loss_history = []
train_accuracy_history = []
val_accuracy_history = []

start_time_train = time.time()

for epoch in range(num_epochs):
    shallow_baseline_model.train() # Set the model to training mode
    running_loss = 0.0
    correct_train = 0
    total_train = 0

    for inputs, labels in train_dataloader:
        # Forward pass
        outputs = shallow_baseline_model(inputs)
        loss = criterion(outputs, labels)

        # Backward and optimize
        optimizer.zero_grad() # Clear gradients
        loss.backward() # Backpropagation
        optimizer.step() # Update weights

        running_loss += loss.item() * inputs.size(0)

        # Calculate training accuracy
        _, predicted = torch.max(outputs.data, 1)
        _, labels_indices = torch.max(labels.data, 1) # Get class indices from original labels
        total_train += labels.size(0)

```

```

    correct_train += (predicted == labels_indices).sum().item()

    epoch_loss = running_loss / len(train_dataloader.dataset)
    epoch_accuracy = 100 * correct_train / total_train
    train_loss_history.append(epoch_loss)
    train_accuracy_history.append(epoch_accuracy)

    # --- Validation Loop ---
    shallow_baseline_model.eval() # Set the model to evaluation mode
    val_running_loss = 0.0
    correct_val = 0
    total_val = 0

    with torch.no_grad(): # Disable gradient calculation for validation
        for inputs_val, labels_val in test_dataloader:
            outputs_val = shallow_baseline_model(inputs_val)
            loss_val = criterion(outputs_val, labels_val)

            val_running_loss += loss_val.item() * inputs_val.size(0)

            # Calculate validation accuracy
            _, predicted_val = torch.max(outputs_val.data, 1)
            _, labels_indices_val = torch.max(labels_val.data, 1)
            total_val += labels_val.size(0)
            correct_val += (predicted_val == labels_indices_val).sum().item()

    epoch_val_loss = val_running_loss / len(test_dataloader.dataset)
    epoch_val_accuracy = 100 * correct_val / total_val
    val_loss_history.append(epoch_val_loss)
    val_accuracy_history.append(epoch_val_accuracy)

    print(f'Epoch [{epoch+1}/{num_epochs}], Loss: {epoch_loss:.4f}, Accuracy: {epoch_accuracy:.2f}')

end_time_train = time.time()
total_training_time = end_time_train - start_time_train

print("\n--- Training Complete ---")
print(f"Total Training Time: {total_training_time:.4f} seconds")

# --- Final Evaluation on Test Data ---
# The last validation accuracy and loss recorded in the loop are the final test !
final_test_loss_shallow_baseline = val_loss_history[-1]
final_test_accuracy_shallow_baseline = val_accuracy_history[-1]

print(f"\nFinal Shallow Baseline Model Performance on Test Data:")
print(f"  Test Loss: {final_test_loss_shallow_baseline:.4f}")
print(f"  Test Accuracy: {final_test_accuracy_shallow_baseline:.2f}%")

# You can also store the loss and accuracy histories for plotting later
shallow_baseline_history = f

```

```
    train_loss_history = [
        'train_loss': train_loss_history,
        'val_loss': val_loss_history,
        'train_accuracy': train_accuracy_history,
        'val_accuracy': val_accuracy_history
    ]
```

Next steps: [Explain error](#)

```
import torch
import torch.nn as nn
import torch.optim as optim
from torch.utils.data import DataLoader # Import DataLoader if not already imported
import time # Import time

# Assuming train_dataloader and test_dataloader are defined from the data preparation
# Assuming X_processed is available to determine input features if needed (though it's not used here)
# Assuming y_train_tensor is available to determine output classes

print("--- Starting Training for Shallow Baseline Neural Network (Combined Step 1) ---")

# --- Define the Shallow Baseline Model architecture ---
# Get input features from the first batch of the training dataloader
if 'train_dataloader' in locals() and next(iter(train_dataloader), None) is not None:
    input_features_soccer = next(iter(train_dataloader))[0].shape[1]
elif 'X_processed' in locals():
    input_features_soccer = X_processed.shape[1]
    print(f"Using X_processed.shape[1] = {input_features_soccer} for input features")
else:
    raise NameError("train_dataloader or X_processed not found. Cannot determine input features")

# Get output classes from the shape of the one-hot encoded target tensor
if 'y_train_tensor' in locals():
    output_classes_soccer = y_train_tensor.shape[1]
else:
    raise NameError("y_train_tensor not found. Cannot determine output classes")
```

```

# Define the Shallow Baseline Model architecture
class ShallowBaselineNN(nn.Module):
    def __init__(self, input_size, hidden_size, num_classes):
        super(ShallowBaselineNN, self).__init__()
        self.fc1 = nn.Linear(input_size, hidden_size) # First hidden layer
        self.relu = nn.ReLU() # ReLU activation
        self.fc2 = nn.Linear(hidden_size, num_classes) # Output layer

    def forward(self, x):
        out = self.fc1(x)
        out = self.relu(out)
        out = self.fc2(out)
        return out # No softmax here, typically applied with the loss function

# Instantiate the model
# Choose a hidden layer size (e.g., 64)
hidden_layer_size = 64
shallow_baseline_model = ShallowBaselineNN(input_features_soccer, hidden_layer_size)

print("Shallow Baseline Neural Network architecture defined.")
print(shallow_baseline_model)

# --- Training Loop ---
# Define Loss Function and Optimizer
# For multi-class classification with one-hot encoded labels, CrossEntropyLoss
criterion = nn.CrossEntropyLoss()
# Adam optimizer is a good default choice
optimizer = optim.Adam(shallow_baseline_model.parameters(), lr=0.001) # Learning rate

num_epochs = 50 # Define the number of training epochs (can be adjusted)
train_loss_history = []
val_loss_history = []
train_accuracy_history = []
val_accuracy_history = []

start_time_train = time.time()

for epoch in range(num_epochs):
    shallow_baseline_model.train() # Set the model to training mode
    running_loss = 0.0
    correct_train = 0
    total_train = 0

    for inputs, labels in train_dataloader:
        # Forward pass
        outputs = shallow_baseline_model(inputs)
        loss = criterion(outputs, labels)

        # Backward and optimize
        optimizer.zero_grad()
        loss.backward()
        optimizer.step()

    # Compute accuracy and loss for validation set
    with torch.no_grad():
        correct_val = 0
        total_val = 0
        for inputs, labels in val_dataloader:
            outputs = shallow_baseline_model(inputs)
            _, predicted = torch.max(outputs, 1)
            total_val += labels.size(0)
            correct_val += (predicted == labels).sum().item()

    train_accuracy_history.append(correct_train / total_train)
    val_accuracy_history.append(correct_val / total_val)

    train_loss_history.append(running_loss / len(train_dataloader))
    val_loss_history.append(loss.item())

    print(f'Epoch {epoch+1}/{num_epochs} | Train Loss: {running_loss / len(train_dataloader)} | Train Acc: {correct_train / total_train} | Val Loss: {loss.item()} | Val Acc: {correct_val / total_val}')

```

```

        optimizer.zero_grad() # Clear gradients
        loss.backward() # Backpropagation
        optimizer.step() # Update weights

        running_loss += loss.item() * inputs.size(0)

        # Calculate training accuracy
        _, predicted = torch.max(outputs.data, 1)
        _, labels_indices = torch.max(labels.data, 1) # Get class indices from
        total_train += labels.size(0)
        correct_train += (predicted == labels_indices).sum().item()

        epoch_loss = running_loss / len(train_dataloader.dataset)
        epoch_accuracy = 100 * correct_train / total_train
        train_loss_history.append(epoch_loss)
        train_accuracy_history.append(epoch_accuracy)

        # --- Validation Loop ---
        shallow_baseline_model.eval() # Set the model to evaluation mode
        val_running_loss = 0.0
        correct_val = 0
        total_val = 0

        with torch.no_grad(): # Disable gradient calculation for validation
            for inputs_val, labels_val in test_dataloader:
                outputs_val = shallow_baseline_model(inputs_val)
                loss_val = criterion(outputs_val, labels_val)

                val_running_loss += loss_val.item() * inputs_val.size(0)

                # Calculate validation accuracy
                _, predicted_val = torch.max(outputs_val.data, 1)
                _, labels_indices_val = torch.max(labels_val.data, 1)
                total_val += labels_val.size(0)
                correct_val += (predicted_val == labels_indices_val).sum().item()

        epoch_val_loss = val_running_loss / len(test_dataloader.dataset)
        epoch_val_accuracy = 100 * correct_val / total_val
        val_loss_history.append(epoch_val_loss)
        val_accuracy_history.append(epoch_val_accuracy)

    print(f'Epoch [{epoch+1}/{num_epochs}], Loss: {epoch_loss:.4f}, Accuracy: {epoch_accuracy:.4f}')

    end_time_train = time.time()
    total_training_time = end_time_train - start_time_train

    print("\n--- Training Complete ---")
    print(f"Total Training Time: {total_training_time:.4f} seconds")

    # --- Final Evaluation on Test Data ---

```

```

    .. FINAL EVALUATION ON TEST DATA

# The last validation accuracy and loss recorded in the loop are the final test
final_test_loss_shallow_baseline = val_loss_history[-1]
final_test_accuracy_shallow_baseline = val_accuracy_history[-1]

print(f"\nFinal Shallow Baseline Model Performance on Test Data:")
print(f"  Test Loss: {final_test_loss_shallow_baseline:.4f}")
print(f"  Test Accuracy: {final_test_accuracy_shallow_baseline:.2f}%")

# You can also store the loss and accuracy histories for plotting later
shallow_baseline_history = {
    'train_loss': train_loss_history,
    'val_loss': val_loss_history,
    'train_accuracy': train_accuracy_history,
    'val_accuracy': val_accuracy_history
}

```

▼ Define Shallow Neural Network (Enhanced) for Soccer Data

Building upon the baseline, this enhanced shallow network will include techniques like Batch Normalization and Dropout, and use the custom Swish activation function to potentially improve training stability and model performance.

```

import torch
import torch.nn as nn

# Assuming input_features_soccer and output_classes_soccer are defined from pre
# Assuming the Swish custom activation function is defined from a previous step

# Define the Shallow Enhanced Model architecture
class ShallowEnhancedNN(nn.Module):
    def __init__(self, input_size, hidden_size, num_classes, dropout_prob=0.5):
        super(ShallowEnhancedNN, self).__init__()
        self.fc1 = nn.Linear(input_size, hidden_size) # First hidden layer
        # self.bn1 = nn.BatchNorm1d(hidden_size) # Original Batch Normalization
        self.ln1 = nn.LayerNorm(hidden_size) # Replace BatchNorm1d with LayerNorm
        self.swish = Swish() # Custom Swish activation
        self.dropout = nn.Dropout(dropout_prob) # Dropout layer
        self.fc2 = nn.Linear(hidden_size, num_classes) # Output layer

    def forward(self, x):
        out = self.fc1(x)
        # out = self.bn1(out) # Original Batch Normalization forward pass
        out = self.ln1(out) # Layer Normalization forward pass
        out = self.swish(out)
        out = self.dropout(out)
        out = self.fc2(out)
        return out # No softmax here, typically applied with the loss function

```

```

# Instantiate the model
# Use the same hidden layer size as the baseline for comparison, e.g., 64
hidden_layer_size = 64
# You can adjust the dropout probability
dropout_probability = 0.5
shallow_enhanced_model = ShallowEnhancedNN(input_features_soccer, hidden_layer_

print("Shallow Enhanced Neural Network architecture defined (using LayerNorm).")
print(shallow_enhanced_model)

```

Train and Evaluate Shallow Neural Network (Enhanced) for Soccer Data

Now that we have defined the architecture for the shallow enhanced model, we need to train it on our processed soccer training data and evaluate its performance on the testing data.

```

import torch
import torch.nn as nn
import torch.optim as optim
from torch.utils.data import DataLoader # Import DataLoader if not already impo
import time # Import time

# Assuming shallow_enhanced_model is defined from the previous step
# Assuming train_dataloader and test_dataloader are defined from the data prep
# Assuming criterion (CrossEntropyLoss) and optimizer (Adam) are defined (or re

print("--- Starting Training for Shallow Enhanced Neural Network ---")

# Define Loss Function and Optimizer (redefine to ensure they are associated wi
criterion_enhanced = nn.CrossEntropyLoss()
optimizer_enhanced = optim.Adam(shallow_enhanced_model.parameters(), lr=0.001)

# --- Training Loop ---
num_epochs = 50 # Define the number of training epochs (can be adjusted)
train_loss_history_enhanced = []
val_loss_history_enhanced = []
train_accuracy_history_enhanced = []
val_accuracy_history_enhanced = []

start_time_train_enhanced = time.time()

for epoch in range(num_epochs):
    shallow_enhanced_model.train() # Set the model to training mode
    running_loss_enhanced = 0.0

```

```

correct_train_enhanced = 0
total_train_enhanced = 0

for i, (inputs, labels) in enumerate(train_dataloader):
    # Check and print batch size
    # print(f"Epoch {epoch+1}, Batch {i+1}: Batch size = {inputs.size(0)}")

    # Forward pass
    outputs = shallow_enhanced_model(inputs)
    loss = criterion_enhanced(outputs, labels)

    # Backward and optimize
    optimizer_enhanced.zero_grad() # Clear gradients
    loss.backward() # Backpropagation
    optimizer_enhanced.step() # Update weights

    running_loss_enhanced += loss.item() * inputs.size(0)

    # Calculate training accuracy
    _, predicted = torch.max(outputs.data, 1)
    _, labels_indices = torch.max(labels.data, 1) # Get class indices from
    total_train_enhanced += labels.size(0)
    correct_train_enhanced += (predicted == labels_indices).sum().item()

epoch_loss_enhanced = running_loss_enhanced / len(train_dataloader.dataset)
epoch_accuracy_enhanced = 100 * correct_train_enhanced / total_train_enhanced
train_loss_history_enhanced.append(epoch_loss_enhanced)
train_accuracy_history_enhanced.append(epoch_accuracy_enhanced)

# --- Validation Loop ---
shallow_enhanced_model.eval() # Set the model to evaluation mode
val_running_loss_enhanced = 0.0
correct_val_enhanced = 0
total_val_enhanced = 0

with torch.no_grad(): # Disable gradient calculation for validation
    for inputs_val, labels_val in test_dataloader:
        outputs_val = shallow_enhanced_model(inputs_val)
        loss_val = criterion_enhanced(outputs_val, labels_val)

        val_running_loss_enhanced += loss_val.item() * inputs_val.size(0)

        # Calculate validation accuracy
        _, predicted_val = torch.max(outputs_val.data, 1)
        _, labels_indices_val = torch.max(labels_val.data, 1)
        total_val_enhanced += labels_val.size(0)
        correct_val_enhanced += (predicted_val == labels_indices_val).sum()

epoch_val_loss_enhanced = val_running_loss_enhanced / len(test_dataloader.dataset)
epoch_val_accuracy_enhanced = 100 * correct_val_enhanced / total_val_enhanced

```

```

epoch_val_accuracy_enhanced = 100 * correct_val_enhanced / total_val_enhanced
val_loss_history_enhanced.append(epoch_val_loss_enhanced)
val_accuracy_history_enhanced.append(epoch_val_accuracy_enhanced)

print(f'Epoch [{epoch+1}/{num_epochs}], Loss: {epoch_loss_enhanced:.4f}, Ac
end_time_train_enhanced = time.time()
total_training_time_enhanced = end_time_train_enhanced - start_time_train_enhan

print("\n--- Training Complete ---")
print(f"Total Training Time: {total_training_time_enhanced:.4f} seconds")

# --- Final Evaluation on Test Data ---
final_test_loss_shallow_enhanced = val_loss_history_enhanced[-1]
final_test_accuracy_shallow_enhanced = val_accuracy_history_enhanced[-1]

print(f"\nFinal Shallow Enhanced Model Performance on Test Data:")
print(f" Test Loss: {final_test_loss_shallow_enhanced:.4f}")
print(f" Test Accuracy: {final_test_accuracy_shallow_enhanced:.2f}%")

# Store history for plotting
shallow_enhanced_history = {
    'train_loss': train_loss_history_enhanced,
    'val_loss': val_loss_history_enhanced,
    'train_accuracy': train_accuracy_history_enhanced,
    'val_accuracy': val_accuracy_history_enhanced
}

```

Plot Shallow Neural Network Training Results

Let's visualize the training and validation loss and accuracy curves for the shallow baseline and enhanced neural networks on the soccer data. This will help us understand their learning progress and identify potential overfitting.

```

import matplotlib.pyplot as plt
import seaborn as sns

# Assuming shallow_baseline_history and shallow_enhanced_history dictionaries e

print(" --- Plotting Shallow Neural Network Training History ---")

# Plot Loss Curves
plt.figure(figsize=(12, 6))
plt.plot(shallow_baseline_history['train_loss'], label='Baseline Train Loss')
plt.plot(shallow_baseline_history['val_loss'], label='Baseline Validation Loss')
plt.plot(shallow_enhanced_history['train_loss'], label='Enhanced Train Loss')

```

```

        plt.plot(shallow_enhanced_history['val_loss'], label='Enhanced Validation Loss')
        plt.title('Shallow Network Loss Curves (Soccer Data)', fontsize=16)
        plt.xlabel('Epoch', fontsize=12)
        plt.ylabel('Loss', fontsize=12)
        plt.legend()
        plt.grid(True, which="both", ls="--")
        plt.show()

    # Plot Accuracy Curves
    plt.figure(figsize=(12, 6))
    plt.plot(shallow_baseline_history['train_accuracy'], label='Baseline Train Accu')
    plt.plot(shallow_baseline_history['val_accuracy'], label='Baseline Validation A')
    plt.plot(shallow_enhanced_history['train_accuracy'], label='Enhanced Train Accu')
    plt.plot(shallow_enhanced_history['val_accuracy'], label='Enhanced Validation A')
    plt.title('Shallow Network Accuracy Curves (Soccer Data)', fontsize=16)
    plt.xlabel('Epoch', fontsize=12)
    plt.ylabel('Accuracy (%)', fontsize=12)
    plt.legend()
    plt.grid(True, which="both", ls="--")
    plt.show()

print("\n--- Plotting Complete! ---")

```

✓ Define Deep Neural Network (Baseline) for Soccer Data

Following the plan, we will now define a **Deep Baseline Neural Network** for the soccer data. This architecture will have at least four hidden layers and use the standard ReLU activation function. It will not include batch normalization or dropout.

```

import torch
import torch.nn as nn

# Assuming input_features_soccer and output_classes_soccer are defined from pre

# Define the Deep Baseline Model architecture
class DeepBaselineNN(nn.Module):
    def __init__(self, input_size, hidden_size_1, hidden_size_2, hidden_size_3,
                 super(DeepBaselineNN, self).__init__()
    self.fc1 = nn.Linear(input_size, hidden_size_1) # First hidden layer
    self.relu1 = nn.ReLU()
    self.fc2 = nn.Linear(hidden_size_1, hidden_size_2) # Second hidden layer
    self.relu2 = nn.ReLU()
    self.fc3 = nn.Linear(hidden_size_2, hidden_size_3) # Third hidden layer
    self.relu3 = nn.ReLU()
    self.fc4 = nn.Linear(hidden_size_3, hidden_size_4) # Fourth hidden layer
    self.relu4 = nn.ReLU()
    # Add more layers if needed to meet the ">= 4 hidden layers" requiremen

```

```

        # but 4 is sufficient for the minimum. Let's add one more for good meas
        self.fc5 = nn.Linear(hidden_size_4, hidden_size_4 // 2) # Fifth hidden
        self.relu5 = nn.ReLU()
        self.fc6 = nn.Linear(hidden_size_4 // 2, num_classes) # Output layer

    def forward(self, x):
        out = self.fc1(x)
        out = self.relu1(out)
        out = self.fc2(out)
        out = self.relu2(out)
        out = self.fc3(out)
        out = self.relu3(out)
        out = self.fc4(out)
        out = self.relu4(out)
        out = self.fc5(out)
        out = self.relu5(out)
        out = self.fc6(out)
        return out # No softmax here

    # Instantiate the model
    # Choose hidden layer sizes (can be adjusted)
    hidden_size_1 = 128
    hidden_size_2 = 64
    hidden_size_3 = 32
    hidden_size_4 = 16

    deep_baseline_model = DeepBaselineNN(input_features_soccer, hidden_size_1, hidde
print("Deep Baseline Neural Network architecture defined.")
print(deep_baseline_model)

```

Train and Evaluate Deep Neural Network (Baseline) for Soccer Data

Now that we have defined the architecture for the deep baseline model, we need to train it on our processed soccer training data and evaluate its performance on the testing data.

```

import torch
import torch.nn as nn
import torch.optim as optim
from torch.utils.data import DataLoader
import time

# Assuming deep_baseline_model is defined from the previous step

```

```

# Assuming train_dataloader and test_dataloader are defined from the data prep
# Assuming criterion (CrossEntropyLoss) and optimizer (Adam) are suitable (or r

print(" --- Starting Training for Deep Baseline Neural Network ---")

# Define Loss Function and Optimizer (redefine to ensure they are associated wi
criterion_deep_baseline = nn.CrossEntropyLoss()
optimizer_deep_baseline = optim.Adam(deep_baseline_model.parameters(), lr=0.001

# --- Training Loop ---
num_epochs = 50 # Define the number of training epochs (can be adjusted)
train_loss_history_deep_baseline = []
val_loss_history_deep_baseline = []
train_accuracy_history_deep_baseline = []
val_accuracy_history_deep_baseline = []

start_time_train_deep_baseline = time.time()

for epoch in range(num_epochs):
    deep_baseline_model.train() # Set the model to training mode
    running_loss_deep_baseline = 0.0
    correct_train_deep_baseline = 0
    total_train_deep_baseline = 0

    for inputs, labels in train_dataloader:
        # Forward pass
        outputs = deep_baseline_model(inputs)
        loss = criterion_deep_baseline(outputs, labels)

        # Backward and optimize
        optimizer_deep_baseline.zero_grad() # Clear gradients
        loss.backward() # Backpropagation
        optimizer_deep_baseline.step() # Update weights

        running_loss_deep_baseline += loss.item() * inputs.size(0)

        # Calculate training accuracy
        _, predicted = torch.max(outputs.data, 1)
        _, labels_indices = torch.max(labels.data, 1) # Get class indices from
        total_train_deep_baseline += labels.size(0)
        correct_train_deep_baseline += (predicted == labels_indices).item()

    epoch_loss_deep_baseline = running_loss_deep_baseline / len(train_dataloader)
    epoch_accuracy_deep_baseline = 100 * correct_train_deep_baseline / total_tr
    train_loss_history_deep_baseline.append(epoch_loss_deep_baseline)
    train_accuracy_history_deep_baseline.append(epoch_accuracy_deep_baseline)

    # --- Validation Loop ---
    deep_baseline_model.eval() # Set the model to evaluation mode
    val_running_loss_deep_baseline = 0.0
    correct_val_deep_baseline = 0

```

```

-----  

total_val_deep_baseline = 0

with torch.no_grad(): # Disable gradient calculation for validation
    for inputs_val, labels_val in test_dataloader:
        outputs_val = deep_baseline_model(inputs_val)
        loss_val = criterion_deep_baseline(outputs_val, labels_val)

        val_running_loss_deep_baseline += loss_val.item() * inputs_val.size(0)

        # Calculate validation accuracy
        _, predicted_val = torch.max(outputs_val.data, 1)
        _, labels_indices_val = torch.max(labels_val.data, 1)
        total_val_deep_baseline += labels_val.size(0)
        correct_val_deep_baseline += (predicted_val == labels_indices_val).sum()

epoch_val_loss_deep_baseline = val_running_loss_deep_baseline / len(test_dataloader)
epoch_val_accuracy_deep_baseline = 100 * correct_val_deep_baseline / total_val_deep_baseline
val_loss_history_deep_baseline.append(epoch_val_loss_deep_baseline)
val_accuracy_history_deep_baseline.append(epoch_val_accuracy_deep_baseline)

print(f'Epoch [{epoch+1}/{num_epochs}], Loss: {epoch_loss_deep_baseline:.4f}')
```

end_time_train_deep_baseline = time.time()
total_training_time_deep_baseline = end_time_train_deep_baseline - start_time_train_deep_baseline

print("\n--- Training Complete ---")
print(f"Total Training Time: {total_training_time_deep_baseline:.4f} seconds")

--- Final Evaluation on Test Data ---
final_test_loss_deep_baseline = val_loss_history_deep_baseline[-1]
final_test_accuracy_deep_baseline = val_accuracy_history_deep_baseline[-1]

print(f"\nFinal Deep Baseline Model Performance on Test Data:")
print(f" Test Loss: {final_test_loss_deep_baseline:.4f}")
print(f" Test Accuracy: {final_test_accuracy_deep_baseline:.2f}%")

You can also store the loss and accuracy histories for plotting later
deep_baseline_history = {
 'train_loss': train_loss_history_deep_baseline,
 'val_loss': val_loss_history_deep_baseline,
 'train_accuracy': train_accuracy_history_deep_baseline,
 'val_accuracy': val_accuracy_history_deep_baseline
}

✓ Define Deep Neural Network (Enhanced) for Soccer Data

Building upon the deep baseline, this enhanced deep network will include Batch

Normalization layers, Dropout layers, and use the custom Swish activation function after each hidden layer (except the output).

```
import torch
import torch.nn as nn

# Assuming input_features_soccer and output_classes_soccer are defined from previous steps
# Assuming the Swish custom activation function is defined from a previous step

# Define the Deep Enhanced Model architecture
class DeepEnhancedNN(nn.Module):
    def __init__(self, input_size, hidden_size_1, hidden_size_2, hidden_size_3,
                 hidden_size_4, num_classes):
        super(DeepEnhancedNN, self).__init__()
        self.fc1 = nn.Linear(input_size, hidden_size_1) # First hidden layer
        self.ln1 = nn.LayerNorm(hidden_size_1) # Layer Normalization (using LayerNorm)
        self.swish1 = Swish() # Custom Swish activation
        self.dropout1 = nn.Dropout(dropout_prob) # Dropout

        self.fc2 = nn.Linear(hidden_size_1, hidden_size_2) # Second hidden layer
        self.ln2 = nn.LayerNorm(hidden_size_2) # Layer Normalization
        self.swish2 = Swish()
        self.dropout2 = nn.Dropout(dropout_prob)

        self.fc3 = nn.Linear(hidden_size_2, hidden_size_3) # Third hidden layer
        self.ln3 = nn.LayerNorm(hidden_size_3) # Layer Normalization
        self.swish3 = Swish()
        self.dropout3 = nn.Dropout(dropout_prob)

        self.fc4 = nn.Linear(hidden_size_3, hidden_size_4) # Fourth hidden layer
        self.ln4 = nn.LayerNorm(hidden_size_4) # Layer Normalization
        self.swish4 = Swish()
        self.dropout4 = nn.Dropout(dropout_prob)

        # Adding the fifth and sixth layers as in the baseline for consistency
        self.fc5 = nn.Linear(hidden_size_4, hidden_size_4 // 2) # Fifth hidden layer
        self.ln5 = nn.LayerNorm(hidden_size_4 // 2) # Layer Normalization
        self.swish5 = Swish()
        self.dropout5 = nn.Dropout(dropout_prob)

        self.fc6 = nn.Linear(hidden_size_4 // 2, num_classes) # Output layer

    def forward(self, x):
        out = self.fc1(x)
        out = self.ln1(out)
        out = self.swish1(out)
        out = self.dropout1(out)

        out = self.fc2(out)
        out = self.ln2(out)
        ... # Add remaining layers
```

```

        out = self.swish2(out)
        out = self.dropout2(out)

        out = self.fc3(out)
        out = self.ln3(out)
        out = self.swish3(out)
        out = self.dropout3(out)

        out = self.fc4(out)
        out = self.ln4(out)
        out = self.swish4(out)
        out = self.dropout4(out)

        out = self.fc5(out)
        out = self.ln5(out)
        out = self.swish5(out)
        out = self.dropout5(out)

        out = self.fc6(out)
        return out # No softmax here

# Instantiate the model with the same hidden layer sizes as the deep baseline
hidden_size_1 = 128
hidden_size_2 = 64
hidden_size_3 = 32
hidden_size_4 = 16
dropout_probability = 0.5

deep_enhanced_model = DeepEnhancedNN(input_features_soccer, hidden_size_1, hidden_size_2, hidden_size_3, hidden_size_4, dropout_probability)

print("Deep Enhanced Neural Network architecture defined (using LayerNorm).")
print(deep_enhanced_model)

```

Train and Evaluate Deep Neural Network (Enhanced) for Soccer Data

Now that we have defined the architecture for the deep enhanced model, we need to train it on our processed soccer training data and evaluate its performance on the testing data.

```

import torch
import torch.nn as nn
import torch.optim as optim
from torch.utils.data import DataLoader
import time

# Assuming deep_enhanced_model is defined from the previous step

```

```

# Assuming deep_enhanced_model is defined from the previous step
# Assuming train_dataloader and test_dataloader are defined from the data preparation
# Assuming criterion (CrossEntropyLoss) and optimizer (Adam) are suitable (or redefined)

print("--- Starting Training for Deep Enhanced Neural Network ---")

# Define Loss Function and Optimizer (redefine to ensure they are associated with the new model)
criterion_deep_enhanced = nn.CrossEntropyLoss()
optimizer_deep_enhanced = optim.Adam(deep_enhanced_model.parameters(), lr=0.001)

# --- Training Loop ---
num_epochs = 50 # Define the number of training epochs (can be adjusted)
train_loss_history_deep_enhanced = []
val_loss_history_deep_enhanced = []
train_accuracy_history_deep_enhanced = []
val_accuracy_history_deep_enhanced = []

start_time_train_deep_enhanced = time.time()

for epoch in range(num_epochs):
    deep_enhanced_model.train() # Set the model to training mode
    running_loss_deep_enhanced = 0.0
    correct_train_deep_enhanced = 0
    total_train_deep_enhanced = 0

    for inputs, labels in train_dataloader:
        # Forward pass
        outputs = deep_enhanced_model(inputs)
        loss = criterion_deep_enhanced(outputs, labels)

        # Backward and optimize
        optimizer_deep_enhanced.zero_grad() # Clear gradients
        loss.backward() # Backpropagation
        optimizer_deep_enhanced.step() # Update weights

        running_loss_deep_enhanced += loss.item() * inputs.size(0)

        # Calculate training accuracy
        _, predicted = torch.max(outputs.data, 1)
        _, labels_indices = torch.max(labels.data, 1) # Get class indices from labels
        total_train_deep_enhanced += labels.size(0)
        correct_train_deep_enhanced += (predicted == labels_indices).sum().item()

    epoch_loss_deep_enhanced = running_loss_deep_enhanced / len(train_dataloader)
    epoch_accuracy_deep_enhanced = 100 * correct_train_deep_enhanced / total_train_deep_enhanced
    train_loss_history_deep_enhanced.append(epoch_loss_deep_enhanced)
    train_accuracy_history_deep_enhanced.append(epoch_accuracy_deep_enhanced)

# --- Validation Loop ---
deep_enhanced_model.eval() # Set the model to evaluation mode
val_running_loss_deep_enhanced = 0.0

```

```

correct_val_deep_enhanced = 0
total_val_deep_enhanced = 0

with torch.no_grad(): # Disable gradient calculation for validation
    for inputs_val, labels_val in test_dataloader:
        outputs_val = deep_enhanced_model(inputs_val)
        loss_val = criterion_deep_enhanced(outputs_val, labels_val)

        val_running_loss_deep_enhanced += loss_val.item() * inputs_val.size(0)

        # Calculate validation accuracy
        _, predicted_val = torch.max(outputs_val.data, 1)
        _, labels_indices_val = torch.max(labels_val.data, 1)
        total_val_deep_enhanced += labels_val.size(0)
        correct_val_deep_enhanced += (predicted_val == labels_indices_val).sum().item()

epoch_val_loss_deep_enhanced = val_running_loss_deep_enhanced / len(test_dataloader)
epoch_val_accuracy_deep_enhanced = 100 * correct_val_deep_enhanced / total_val_deep_enhanced
val_loss_history_deep_enhanced.append(epoch_val_loss_deep_enhanced)
val_accuracy_history_deep_enhanced.append(epoch_val_accuracy_deep_enhanced)

print(f'Epoch [{epoch+1}/{num_epochs}], Loss: {epoch_loss_deep_enhanced:.4f}')
print(f'Val Loss: {epoch_val_loss_deep_enhanced:.4f}')

end_time_train_deep_enhanced = time.time()
total_training_time_deep_enhanced = end_time_train_deep_enhanced - start_time_train_deep_enhanced

print("\n--- Training Complete ---")
print(f"Total Training Time: {total_training_time_deep_enhanced:.4f} seconds")

# --- Final Evaluation on Test Data ---
final_test_loss_deep_enhanced = val_loss_history_deep_enhanced[-1]
final_test_accuracy_deep_enhanced = val_accuracy_history_deep_enhanced[-1]

print("\nFinal Deep Enhanced Model Performance on Test Data:")
print(f"Test Loss: {final_test_loss_deep_enhanced:.4f}")
print(f"Test Accuracy: {final_test_accuracy_deep_enhanced:.2f}%")

# Store history for plotting
deep_enhanced_history = {
    'train_loss': train_loss_history_deep_enhanced,
    'val_loss': val_loss_history_deep_enhanced,
    'train_accuracy': train_accuracy_history_deep_enhanced,
    'val_accuracy': val_accuracy_history_deep_enhanced
}

```

▼ Plot Deep Neural Network Training Results

Let's visualize the training and validation loss and accuracy curves for the deep baseline

and enhanced neural networks on the soccer data. This will help us compare their learning progress and identify potential overfitting, similar to the shallow networks.

```
import matplotlib.pyplot as plt
import seaborn as sns

# Assuming deep_baseline_history and deep_enhanced_history dictionaries exist

print("--- Plotting Deep Neural Network Training History ---")

# Plot Loss Curves
plt.figure(figsize=(12, 6))
plt.plot(deep_baseline_history['train_loss'], label='Baseline Train Loss')
plt.plot(deep_baseline_history['val_loss'], label='Baseline Validation Loss')
plt.plot(deep_enhanced_history['train_loss'], label='Enhanced Train Loss')
plt.plot(deep_enhanced_history['val_loss'], label='Enhanced Validation Loss')
plt.title('Deep Network Loss Curves (Soccer Data)', fontsize=16)
plt.xlabel('Epoch', fontsize=12)
plt.ylabel('Loss', fontsize=12)
plt.legend()
plt.grid(True, which="both", ls="--")
plt.show()

# Plot Accuracy Curves
plt.figure(figsize=(12, 6))
plt.plot(deep_baseline_history['train_accuracy'], label='Baseline Train Accuracy')
plt.plot(deep_baseline_history['val_accuracy'], label='Baseline Validation Accuracy')
plt.plot(deep_enhanced_history['train_accuracy'], label='Enhanced Train Accuracy')
plt.plot(deep_enhanced_history['val_accuracy'], label='Enhanced Validation Accuracy')
plt.title('Deep Network Accuracy Curves (Soccer Data)', fontsize=16)
plt.xlabel('Epoch', fontsize=12)
plt.ylabel('Accuracy (%)', fontsize=12)
plt.legend()
plt.grid(True, which="both", ls="--")
plt.show()

print("\n--- Plotting Complete! ---")
```

❖ Define Wide Neural Network (Baseline) for Soccer Data

Following the plan, we will now define a **Wide Baseline Neural Network** for the soccer data. This architecture will have a single hidden layer with a large number of neurons and use the standard ReLU activation function. It will not include batch normalization or dropout.

```

import torch
import torch.nn as nn

# Assuming input_features_soccer and output_classes_soccer are defined from previous steps

# Define the Wide Baseline Model architecture
class WideBaselineNN(nn.Module):
    def __init__(self, input_size, wide_hidden_size, num_classes):
        super(WideBaselineNN, self).__init__()
        self.fc1 = nn.Linear(input_size, wide_hidden_size) # Single wide hidden layer
        self.relu = nn.ReLU() # ReLU activation
        self.fc2 = nn.Linear(wide_hidden_size, num_classes) # Output layer

    def forward(self, x):
        out = self.fc1(x)
        out = self.relu(out)
        out = self.fc2(out)
        return out # No softmax here

# Instantiate the model
# Choose a large hidden layer size (e.g., 512 or 1024, significantly larger than input size)
wide_hidden_layer_size = 512 # This can be adjusted
wide_baseline_model = WideBaselineNN(input_features_soccer, wide_hidden_layer_size)

print("Wide Baseline Neural Network architecture defined.")
print(wide_baseline_model)

```

Train and Evaluate Wide Neural Network (Baseline) for Soccer Data

Now that we have defined the architecture for the wide baseline model, we need to train it on our processed soccer training data and evaluate its performance on the testing data.

```

import torch
import torch.nn as nn
import torch.optim as optim
from torch.utils.data import DataLoader
import time

# Assuming wide_baseline_model is defined from the previous step
# Assuming train_dataloader and test_dataloader are defined from the data preparation step
# Assuming criterion (CrossEntropyLoss) and optimizer (Adam) are suitable (or redefined)

print("--- Starting Training for Wide Baseline Neural Network ---")

# Define loss function and optimizer (redefined to ensure they are associated with the new model)

```

```

# Define loss function and optimizer (redefining to ensure they are associated with the module)
criterion_wide_baseline = nn.CrossEntropyLoss()
optimizer_wide_baseline = optim.Adam(wide_baseline_model.parameters(), lr=0.001)

# --- Training Loop ---
num_epochs = 50 # Define the number of training epochs (can be adjusted)
train_loss_history_wide_baseline = []
val_loss_history_wide_baseline = []
train_accuracy_history_wide_baseline = []
val_accuracy_history_wide_baseline = []

start_time_train_wide_baseline = time.time()

for epoch in range(num_epochs):
    wide_baseline_model.train() # Set the model to training mode
    running_loss_wide_baseline = 0.0
    correct_train_wide_baseline = 0
    total_train_wide_baseline = 0

    for inputs, labels in train_dataloader:
        # Forward pass
        outputs = wide_baseline_model(inputs)
        loss = criterion_wide_baseline(outputs, labels)

        # Backward and optimize
        optimizer_wide_baseline.zero_grad() # Clear gradients
        loss.backward() # Backpropagation
        optimizer_wide_baseline.step() # Update weights

        running_loss_wide_baseline += loss.item() * inputs.size(0)

        # Calculate training accuracy
        _, predicted = torch.max(outputs.data, 1)
        _, labels_indices = torch.max(labels.data, 1) # Get class indices from labels
        total_train_wide_baseline += labels.size(0)
        correct_train_wide_baseline += (predicted == labels_indices).item()

    epoch_loss_wide_baseline = running_loss_wide_baseline / len(train_dataloader)
    epoch_accuracy_wide_baseline = 100 * correct_train_wide_baseline / total_train_wide_baseline
    train_loss_history_wide_baseline.append(epoch_loss_wide_baseline)
    train_accuracy_history_wide_baseline.append(epoch_accuracy_wide_baseline)

# --- Validation Loop ---
wide_baseline_model.eval() # Set the model to evaluation mode
val_running_loss_wide_baseline = 0.0
correct_val_wide_baseline = 0
total_val_wide_baseline = 0

with torch.no_grad(): # Disable gradient calculation for validation
    for inputs_val, labels_val in test_dataloader:
        outputs_val = wide_baseline_model(inputs_val)

```

```

loss_val = criterion_wide_baseline(outputs_val, labels_val)

val_running_loss_wide_baseline += loss_val.item() * inputs_val.size

# Calculate validation accuracy
_, predicted_val = torch.max(outputs_val.data, 1)
_, labels_indices_val = torch.max(labels_val.data, 1)
total_val_wide_baseline += labels_val.size(0)
correct_val_wide_baseline += (predicted_val == labels_indices_val).

epoch_val_loss_wide_baseline = val_running_loss_wide_baseline / len(test_da
epoch_val_accuracy_wide_baseline = 100 * correct_val_wide_baseline / total_
val_loss_history_wide_baseline.append(epoch_val_loss_wide_baseline)
val_accuracy_history_wide_baseline.append(epoch_val_accuracy_wide_baseline)

print(f'Epoch [{epoch+1}/{num_epochs}], Loss: {epoch_loss_wide_baseline:.4f}

end_time_train_wide_baseline = time.time()
total_training_time_wide_baseline = end_time_train_wide_baseline - start_time_t

print("\n--- Training Complete ---")
print(f"Total Training Time: {total_training_time_wide_baseline:.4f} seconds")

# --- Final Evaluation on Test Data ---
final_test_loss_wide_baseline = val_loss_history_wide_baseline[-1]
final_test_accuracy_wide_baseline = val_accuracy_history_wide_baseline[-1]

print(f"\nFinal Wide Baseline Model Performance on Test Data:")
print(f" Test Loss: {final_test_loss_wide_baseline:.4f}")
print(f" Test Accuracy: {final_test_accuracy_wide_baseline:.2f}%")

# You can also store the loss and accuracy histories for plotting later
wide_baseline_history = {
    'train_loss': train_loss_history_wide_baseline,
    'val_loss': val_loss_history_wide_baseline,
    'train_accuracy': train_accuracy_history_wide_baseline,
    'val_accuracy': val_accuracy_history_wide_baseline
}

```

```

==== STEP 5: RECURRENT NEURAL NETWORK EXTENSION (LSTM / GRU) ====
=====
""")

import pandas as pd
import numpy as np
import torch
import torch.nn as nn
from torch.utils.data import DataLoader, TensorDataset
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.model_selection import train_test_split

# =====
# === 1. CONFIGURATION =====
# =====
SEQ_LEN = 5 # number of past matches per sequence
TARGET_COL = 'winner' # classification target
BATCH_SIZE = 32
EPOCHS = 20
LR = 1e-3
HIDDEN_SIZE = 64
NUM_LAYERS = 1
USE_GRU = False # set True to use GRU instead of LSTM

# =====
# === 2. FEATURE SELECTION & ENCODING =====
# =====
print("\n[INFO] Preparing features and target...")

# Features to include in sequences
if 'total_goals' not in cleaned_df.columns:
    cleaned_df['total_goals'] = cleaned_df['home_score'] + cleaned_df['away_scc']

feature_cols = [
    'goal_difference', 'total_goals', 'expected_goal_diff',
    'home_avg_goal_diff_last5', 'away_avg_goal_diff_last5',
    'match_hour', 'match_day_of_week', 'match_month'
]

# Scale numeric features
scaler = StandardScaler()
cleaned_df[feature_cols] = scaler.fit_transform(cleaned_df[feature_cols])

# One-hot encode target
ohe_target = OneHotEncoder(sparse_output=False, handle_unknown='ignore')
y_all = ohe_target.fit_transform(cleaned_df[[TARGET_COL]])

# =====
# === 3. BUILD SEQUENCES PER TEAM =====
# =====

```

```

print("\n[INFO] Building sequences...")

X_seq, y_seq = [], []
cleaned_df = cleaned_df.sort_values(by='utcDate')

for team, group in cleaned_df.groupby('homeTeam_name'):
    group = group.sort_values('utcDate')
    feats = group[feature_cols].values
    targets = y_all[group.index]
    for i in range(len(group) - SEQ_LEN):
        X_seq.append(feats[i:i+SEQ_LEN])
        y_seq.append(targets[i+SEQ_LEN])

X_seq = np.array(X_seq, dtype=np.float32)
y_seq = np.array(y_seq, dtype=np.float32)

print(f"[INFO] Sequence tensor shape: {X_seq.shape}, Target tensor shape: {y_seq.shape}")
import numpy as np

print("Any NaNs in X_seq?", np.isnan(X_seq).any())
print("Any infs in X_seq?", np.isinf(X_seq).any())
print("Any NaNs in y_seq?", np.isnan(y_seq).any())
print("Any infs in y_seq?", np.isinf(y_seq).any())

# =====
# === 4. TRAIN/TEST SPLIT =====
# =====
X_train, X_test, y_train, y_test = train_test_split(
    X_seq, y_seq, test_size=0.2, shuffle=True, random_state=42
)

train_dataset = TensorDataset(torch.tensor(X_train), torch.tensor(y_train))
test_dataset = TensorDataset(torch.tensor(X_test), torch.tensor(y_test))

train_loader = DataLoader(train_dataset, batch_size=BATCH_SIZE, shuffle=True)
test_loader = DataLoader(test_dataset, batch_size=BATCH_SIZE)

# =====
# === 5. DEFINE LSTM/GRU MODEL =====
# =====
print("\n[INFO] Initializing model...")

class RNNClassifier(nn.Module):
    def __init__(self, input_size, hidden_size, num_layers, output_size, use_gru):
        super().__init__()
        if use_gru:
            self.rnn = nn.GRU(input_size, hidden_size, num_layers, batch_first=True)
        else:
            self.rnn = nn.LSTM(input_size, hidden_size, num_layers, batch_first=True)
        self.fc = nn.Linear(hidden_size, output_size)

```

```

def forward(self, x):
    out, _ = self.rnn(x)
    out = out[:, -1, :]
    return self.fc(out)

model = RNNClassifier(
    input_size=len(feature_cols),
    hidden_size=HIDDEN_SIZE,
    num_layers=NUM_LAYERS,
    output_size=y_seq.shape[1],
    use_gru=USE_GRU
)

# =====
# === 6. TRAINING LOOP =====
# =====
criterion = nn.CrossEntropyLoss()
optimizer = torch.optim.Adam(model.parameters(), lr=LR)

print("\n[INFO] Starting training...")
for epoch in range(EPOCHS):
    model.train()
    running_loss, correct, total = 0.0, 0, 0
    for X_batch, y_batch in train_loader:
        optimizer.zero_grad()
        outputs = model(X_batch)
        loss = criterion(outputs, torch.argmax(y_batch, dim=1))
        loss.backward()
        optimizer.step()
        running_loss += loss.item() * X_batch.size(0)
        preds = torch.argmax(outputs, dim=1)
        correct += (preds == torch.argmax(y_batch, dim=1)).sum().item()
        total += y_batch.size(0)
    train_loss = running_loss / total
    train_acc = correct / total

    # Validation
    model.eval()
    val_loss, val_correct, val_total = 0.0, 0, 0
    with torch.no_grad():
        for X_batch, y_batch in test_loader:
            outputs = model(X_batch)
            loss = criterion(outputs, torch.argmax(y_batch, dim=1))
            val_loss += loss.item() * X_batch.size(0)
            preds = torch.argmax(outputs, dim=1)
            val_correct += (preds == torch.argmax(y_batch, dim=1)).sum().item()
            val_total += y_batch.size(0)
    val_loss /= val_total
    val_acc = val_correct / val_total

```

```

        print(f"Epoch {epoch+1}/{EPOCHS} | Train Loss: {train_loss:.4f} Acc: {train_acc:.4f}\n"
              f"Val Loss: {val_loss:.4f} Acc: {val_acc:.4f}")

# =====
# === 7. FINAL EVALUATION =====
# =====
print("\n[RESULTS] Final Validation Accuracy:", round(val_acc, 4))
print("[RESULTS] Final Validation Loss:", round(val_loss, 4))
print("\n[INFO] Step 5 complete – ready to compare against MLPs and kernel regr

```

```

=====
== STEP 5: RECURRENT NEURAL NETWORK EXTENSION (LSTM / GRU) ==
=====
```

[INFO] Preparing features and target...

[INFO] Building sequences...

[INFO] Sequence tensor shape: (5602, 5, 8), Target tensor shape: (5602, 3)
Any NaNs in X_seq? True
Any inf in X_seq? False
Any NaNs in y_seq? False
Any inf in y_seq? False

[INFO] Initializing model...

[INFO] Starting training...

Epoch 1/20	Train Loss: nan	Acc: 0.3053	Val Loss: nan	Acc: 0.3310
Epoch 2/20	Train Loss: nan	Acc: 0.3053	Val Loss: nan	Acc: 0.3310
Epoch 3/20	Train Loss: nan	Acc: 0.3053	Val Loss: nan	Acc: 0.3310
Epoch 4/20	Train Loss: nan	Acc: 0.3053	Val Loss: nan	Acc: 0.3310
Epoch 5/20	Train Loss: nan	Acc: 0.3053	Val Loss: nan	Acc: 0.3310
Epoch 6/20	Train Loss: nan	Acc: 0.3053	Val Loss: nan	Acc: 0.3310
Epoch 7/20	Train Loss: nan	Acc: 0.3053	Val Loss: nan	Acc: 0.3310
Epoch 8/20	Train Loss: nan	Acc: 0.3053	Val Loss: nan	Acc: 0.3310
Epoch 9/20	Train Loss: nan	Acc: 0.3053	Val Loss: nan	Acc: 0.3310
Epoch 10/20	Train Loss: nan	Acc: 0.3053	Val Loss: nan	Acc: 0.3310
Epoch 11/20	Train Loss: nan	Acc: 0.3053	Val Loss: nan	Acc: 0.3310
Epoch 12/20	Train Loss: nan	Acc: 0.3053	Val Loss: nan	Acc: 0.3310
Epoch 13/20	Train Loss: nan	Acc: 0.3053	Val Loss: nan	Acc: 0.3310
Epoch 14/20	Train Loss: nan	Acc: 0.3053	Val Loss: nan	Acc: 0.3310
Epoch 15/20	Train Loss: nan	Acc: 0.3053	Val Loss: nan	Acc: 0.3310
Epoch 16/20	Train Loss: nan	Acc: 0.3053	Val Loss: nan	Acc: 0.3310
Epoch 17/20	Train Loss: nan	Acc: 0.3053	Val Loss: nan	Acc: 0.3310
Epoch 18/20	Train Loss: nan	Acc: 0.3053	Val Loss: nan	Acc: 0.3310
Epoch 19/20	Train Loss: nan	Acc: 0.3053	Val Loss: nan	Acc: 0.3310
Epoch 20/20	Train Loss: nan	Acc: 0.3053	Val Loss: nan	Acc: 0.3310

[RESULTS] Final Validation Accuracy: 0.331

[RESULTS] Final Validation Loss: nan

[INFO] Step 5 complete – ready to compare against MLPs and kernel regressors.

Hmm.. after a couple of definition errors, it still has a pretty blah accuracy number. Let's try a few tweaks - such as the dropout and weight decay - to see if we can enhance and stabilize a more robust version.

```
# =====
# === STEP 5: RECURRENT NN EXTENSION (LSTM / GRU) - WITH CM ====
# =====

import numpy as np
import torch
import torch.nn as nn
from torch.utils.data import DataLoader, TensorDataset
from sklearn.preprocessing import StandardScaler
from torch.optim.lr_scheduler import ReduceLROnPlateau
from sklearn.metrics import confusion_matrix, classification_report
import pandas as pd

print("\n[INFO] Preparing features and target...")

# -----
# 1. Define features & target
# -----
feature_cols = [
    'goal_difference', 'total_goals', 'expected_goal_diff',
    'home_avg_goal_diff_last5', 'away_avg_goal_diff_last5',
    'match_hour', 'match_day_of_week', 'match_month'
]
target_classes = ['HomeWin', 'Draw', 'AwayWin']

# Ensure all required features exist
for col in feature_cols:
    if col not in cleaned_df.columns:
        raise ValueError(f"Missing required feature: {col}")

# Fill NaNs
cleaned_df[feature_cols] = cleaned_df[feature_cols].fillna(0)

# Scale numeric features
scaler = StandardScaler()
cleaned_df[feature_cols] = scaler.fit_transform(cleaned_df[feature_cols])

print("[DEBUG] Feature sample after scaling:")
print(cleaned_df[feature_cols].head())

# -----
# 2. Build sequences
# -----
```

```

seq_length = 5

def build_sequences(df, features, target, seq_len):
    X, y = [], []
    for i in range(len(df) - seq_len):
        X.append(df[features].iloc[i:i+seq_len].values)
        y.append(df[target].iloc[i+seq_len])
    return np.array(X), np.array(y)

# Map winner to class index
outcome_map = {'HOME_TEAM': 0, 'DRAW': 1, 'AWAY_TEAM': 2}
cleaned_df['target_class'] = cleaned_df['winner'].map(outcome_map).fillna(1).as

X_seq, y_seq = build_sequences(cleaned_df, feature_cols, 'target_class', seq_le

print(f"[INFO] Sequence tensor shape: {X_seq.shape}, Target tensor shape: {y_se
print("Any NaNs in X_seq?", np.isnan(X_seq).any())
print("Any infs in X_seq?", np.isinf(X_seq).any())

# -----
# 3. Convert to PyTorch tensors
# -----
X_tensor = torch.tensor(X_seq, dtype=torch.float32)
y_tensor = torch.tensor(y_seq, dtype=torch.long)

dataset = TensorDataset(X_tensor, y_tensor)
train_size = int(0.8 * len(dataset))
val_size = len(dataset) - train_size
train_dataset, val_dataset = torch.utils.data.random_split(dataset, [train_size, v

train_loader = DataLoader(train_dataset, batch_size=32, shuffle=True)
val_loader = DataLoader(val_dataset, batch_size=32)

# -----
# 4. Define RNN model
# -----
class MatchRNN(nn.Module):
    def __init__(self, input_size, hidden_size, num_layers, num_classes, rnn_ty
        super(MatchRNN, self).__init__()
        if rnn_type == 'LSTM':
            self.rnn = nn.LSTM(input_size, hidden_size, num_layers, batch_first=t
        else:
            self.rnn = nn.GRU(input_size, hidden_size, num_layers, batch_first=t
        self.dropout = nn.Dropout(dropout)
        self.fc = nn.Linear(hidden_size, num_classes)

    def forward(self, x):
        out, _ = self.rnn(x)
        out = out[:, -1, :]
        out = self.dropout(out)
        out = self.fc(out)


```

```

        out = self.rnn(out)
        return out

input_size = len(feature_cols)
hidden_size = 64
num_layers = 2 # enables recurrent dropout
num_classes = len(target_classes)

model = MatchRNN(input_size, hidden_size, num_layers, num_classes, rnn_type='LS'
criterion = nn.CrossEntropyLoss()
optimizer = torch.optim.Adam(model.parameters(), lr=1e-3, weight_decay=1e-5)
scheduler = ReduceLROnPlateau(optimizer, mode='min', patience=3, factor=0.5)

print("\n[INFO] Starting training...")

# -----
# 5. Training loop with early stopping & per-class accuracy
# -----
best_val_loss = float('inf')
patience = 5
patience_counter = 0

for epoch in range(1, 21):
    model.train()
    train_loss, correct, total = 0, 0, 0
    for X_batch, y_batch in train_loader:
        optimizer.zero_grad()
        outputs = model(X_batch)
        loss = criterion(outputs, y_batch)
        loss.backward()
        optimizer.step()

        train_loss += loss.item() * X_batch.size(0)
        _, predicted = torch.max(outputs, 1)
        total += y_batch.size(0)
        correct += (predicted == y_batch).sum().item()

    train_acc = correct / total
    train_loss /= total

    # Validation
    model.eval()
    val_loss, val_correct, val_total = 0, 0, 0
    class_correct = [0] * num_classes
    class_total = [0] * num_classes
    all_preds, all_labels = [], []

    with torch.no_grad():
        for X_batch, y_batch in val_loader:
            outputs = model(X_batch)
            loss = criterion(outputs, y_batch)

            val_loss += loss.item() * X_batch.size(0)
            _, predicted = torch.max(outputs, 1)
            total += y_batch.size(0)
            correct += (predicted == y_batch).sum().item()

            all_preds.append(predicted)
            all_labels.append(y_batch)

    val_acc = correct / total
    val_loss /= total

    if val_loss < best_val_loss:
        best_val_loss = val_loss
        patience_counter = 0
    else:
        patience_counter += 1

    if patience_counter == patience:
        print(f"\n[INFO] Early stopping at epoch {epoch} due to no improvement in validation loss for {patience} epochs")
        break

```

```

        val_loss += loss.item() * X_batch.size(0)
        _, predicted = torch.max(outputs, 1)
        val_total += y_batch.size(0)
        val_correct += (predicted == y_batch).sum().item()

        all_preds.extend(predicted.cpu().numpy())
        all_labels.extend(y_batch.cpu().numpy())

    for i in range(len(y_batch)):
        label = y_batch[i].item()
        class_total[label] += 1
        if predicted[i].item() == label:
            class_correct[label] += 1

    val_acc = val_correct / val_total
    val_loss /= val_total

    # Per-class accuracy
    per_class_acc = {target_classes[i]: (class_correct[i] / class_total[i]) if c
                     for i in range(num_classes)}

    print(f"Epoch {epoch}/20 | Train Loss: {train_loss:.4f} Acc: {train_acc:.4f}
          f"Val Loss: {val_loss:.4f} Acc: {val_acc:.4f} | Per-class Acc: {per_c

    # Manual LR change print
    old_lr = optimizer.param_groups[0]['lr']
    scheduler.step(val_loss)
    new_lr = optimizer.param_groups[0]['lr']
    if new_lr != old_lr:
        print(f"[INFO] Learning rate reduced from {old_lr} to {new_lr}")

    # Early stopping
    if val_loss < best_val_loss:
        best_val_loss = val_loss
        patience_counter = 0
        best_preds = list(all_preds)
        best_labels = list(all_labels)
    else:
        patience_counter += 1
        if patience_counter >= patience:
            print("[INFO] Early stopping triggered.")
            break

# -----
# 6. Confusion Matrix & Report
# -----
cm = confusion_matrix(best_labels, best_preds)
cm_df = pd.DataFrame(cm, index=target_classes, columns=target_classes)

print("\n[CONFUSION MATRIX]")
print(cm_df)

```

```

print("\n[CLASSIFICATION REPORT]")
print(classification_report(best_labels, best_preds, target_names=target_classes))

print("\n[RESULTS] Final Validation Accuracy:", val_acc)
print("[RESULTS] Final Validation Loss:", val_loss)

```

[INFO] Preparing features and target...

[DEBUG] Feature sample after scaling:

	goal_difference	total_goals	expected_goal_diff	\
1224	0.397325	0.137328	5.366714e-17	
1225	-1.820038	0.137328	5.366714e-17	
1226	-0.157016	-0.467099	5.366714e-17	
1227	0.951666	-0.467099	5.366714e-17	
1228	0.397325	0.137328	5.366714e-17	

	home_avg_goal_diff_last5	away_avg_goal_diff_last5	match_hour	\
1224	4.043862e-17	-1.530110e-17	0.788423	
1225	4.043862e-17	-1.530110e-17	0.788423	
1226	4.043862e-17	-1.530110e-17	1.321930	
1227	4.043862e-17	-1.530110e-17	1.321930	
1228	4.043862e-17	-1.530110e-17	1.321930	

	match_day_of_week	match_month
1224	0.216728	-0.704832
1225	0.216728	-0.704832
1226	0.216728	-0.704832
1227	0.216728	-0.704832
1228	0.216728	-0.704832

[INFO] Sequence tensor shape: (6597, 5, 8), Target tensor shape: (6597,)

Any NaNs in X_seq? False

Any infs in X_seq? False

[INFO] Starting training...

Epoch	Train Loss	Acc	Val Loss	Acc	Per
1/20	1.0771	0.4252	1.0822	0.4250	Per
2/20	1.0728	0.4402	1.0795	0.4250	Per
3/20	1.0710	0.4402	1.0823	0.4250	Per
4/20	1.0707	0.4410	1.0824	0.4250	Per
5/20	1.0692	0.4395	1.0805	0.4242	Per
6/20	1.0688	0.4402	1.0823	0.4197	Per

[INFO] Learning rate reduced from 0.001 to 0.0005

Epoch	Train Loss	Acc	Val Loss	Acc	Per
7/20	1.0664	0.4415	1.0846	0.4197	Per

[INFO] Early stopping triggered.

[CONFUSION MATRIX]

	HomeWin	Draw	AwayWin
HomeWin	561	0	0
Draw	353	0	0
AwayWin	406	0	0

[CLASSIFICATION REPORT]

precision	recall	f1-score	support
-----------	--------	----------	---------

HomeWin	0.42	1.00	0.60	561
Draw	0.00	0.00	0.00	353
AwayWin	0.00	0.00	0.00	406
			0.42	1320
accuracy			0.42	1320
macro avg	0.14	0.33	0.20	1320
weighted avg	0.18	0.42	0.25	1320

[RESULTS] Final Validation Accuracy: 0.4196969696969697

[RESULTS] Final Validation Loss: 1.0845934651114724

The RNN seems to recall 100% of HomeWin, which is basically just flipping a coin. It's not really predicting anything - even worse than a coin flip - even while using an enhanced RNN architecture.

```
#####
# #####
# ###
# ### STEP 6: FEATURE-TRANSFER EXPERIMENT #####
# ###

print(""""

=====
== STEP 6: FEATURE-TRANSFER EXPERIMENT ==
=====

""")

import pandas as pd
import numpy as np
import torch
import torch.nn as nn
from torch.utils.data import DataLoader, TensorDataset
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score

# =====
# === 1. CONFIGURATION ===
# =====
BATCH_SIZE = 32
EPOCHS = 20
LR = 1e-3
HIDDEN_SIZE = 64

# =====
# === 2. PREPARE FLAT FEATURES & TARGET ===
=====
```

```

# =====
print("\n[INFO] Preparing flat features for MLP training...")

# Choose features (can include engineered ones from Step 5)
flat_feature_cols = [
    'goal_difference', 'total_goals', 'expected_goal_diff',
    'home_avg_goal_diff_last5', 'away_avg_goal_diff_last5',
    'match_hour', 'match_day_of_week', 'match_month'
]

# Scale numeric features
scaler = StandardScaler()
X_flat = scaler.fit_transform(cleaned_df[flat_feature_cols])

# One-hot encode target
ohe_target = OneHotEncoder(sparse_output=False, handle_unknown='ignore')
y_flat = ohe_target.fit_transform(cleaned_df[['winner']])

# Train/test split
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(
    X_flat, y_flat, test_size=0.2, random_state=42, shuffle=True
)

train_dataset = TensorDataset(torch.tensor(X_train, dtype=torch.float32),
                             torch.tensor(y_train, dtype=torch.float32))
test_dataset = TensorDataset(torch.tensor(X_test, dtype=torch.float32),
                           torch.tensor(y_test, dtype=torch.float32))

train_loader = DataLoader(train_dataset, batch_size=BATCH_SIZE, shuffle=True)
test_loader = DataLoader(test_dataset, batch_size=BATCH_SIZE)

# =====
# === 3. DEFINE FEEDFORWARD NETWORK =====
# =====
print("\n[INFO] Initializing feedforward network...")

class FeedforwardNet(nn.Module):
    def __init__(self, input_dim, hidden_dim, output_dim):
        super().__init__()
        self.feature_extractor = nn.Sequential(
            nn.Linear(input_dim, hidden_dim),
            nn.ReLU(),
            nn.Linear(hidden_dim, hidden_dim),
            nn.ReLU()
        )
        self.classifier = nn.Linear(hidden_dim, output_dim)

    def forward(self, x):
        feats = self.feature_extractor(x)
        ...

```

```

        out = self.classifier(feats)
        return out, feats

model = FeedforwardNet(input_dim=X_train.shape[1],
                      hidden_dim=HIDDEN_SIZE,
                      output_dim=y_train.shape[1])

# =====
# === 4. TRAIN MLP ON RAW FEATURES =====
# =====
criterion = nn.CrossEntropyLoss()
optimizer = torch.optim.Adam(model.parameters(), lr=LR)

print("\n[INFO] Training MLP on raw features...")
for epoch in range(EPOCHS):
    model.train()
    running_loss, correct, total = 0.0, 0, 0
    for X_batch, y_batch in train_loader:
        optimizer.zero_grad()
        outputs, _ = model(X_batch)
        loss = criterion(outputs, torch.argmax(y_batch, dim=1))
        loss.backward()
        optimizer.step()
        running_loss += loss.item() * X_batch.size(0)
        preds = torch.argmax(outputs, dim=1)
        correct += (preds == torch.argmax(y_batch, dim=1)).sum().item()
        total += y_batch.size(0)
    train_loss = running_loss / total
    train_acc = correct / total
    print(f"Epoch {epoch+1}/{EPOCHS} | Train Loss: {train_loss:.4f} Acc: {train_acc:.4f}")

# =====
# === 5. FREEZE PENULTIMATE LAYER & EXTRACT EMBEDDINGS =====
# =====
print("\n[INFO] Freezing feature extractor and generating embeddings...")

for param in model.feature_extractor.parameters():
    param.requires_grad = False

with torch.no_grad():
    train_embeddings = model.feature_extractor(torch.tensor(X_train, dtype=torch.float32))
    test_embeddings = model.feature_extractor(torch.tensor(X_test, dtype=torch.float32))

print(f"[INFO] Train embeddings shape: {train_embeddings.shape}")
print(f"[INFO] Test embeddings shape: {test_embeddings.shape}")

# =====
# === 6. TRAIN LINEAR MODEL ON EMBEDDINGS =====
# =====
print("\n[INFO] Training Logistic Regression on learned embeddings...")

```

```

log_reg = LogisticRegression(max_iter=1000)
log_reg.fit(train_embeddings, np.argmax(y_train, axis=1))

y_pred = log_reg.predict(test_embeddings)
transfer_acc = accuracy_score(np.argmax(y_test, axis=1), y_pred)

print(f"[RESULTS] Accuracy of Logistic Regression on learned embeddings: {transfer_acc:.4f}\n")

# =====
# === 7. BASELINE: LINEAR MODEL ON RAW FEATURES =====
# =====
print("\n[INFO] Training Logistic Regression on raw features...")

log_reg_baseline = LogisticRegression(max_iter=1000)
log_reg_baseline.fit(X_train, np.argmax(y_train, axis=1))

y_pred_base = log_reg_baseline.predict(X_test)
baseline_acc = accuracy_score(np.argmax(y_test, axis=1), y_pred_base)

print(f"[RESULTS] Accuracy of Logistic Regression on raw features: {baseline_acc:.4f}\n")

# =====
# === 8. COMPARISON OUTPUT =====
# =====
print("=====")
print(" FEATURE-TRANSFER COMPARISON RESULTS")
print("-----")
print(f" Logistic Regression on raw features: {baseline_acc:.4f}")
print(f" Logistic Regression on learned features: {transfer_acc:.4f}")
print("=====")

```

```

=====
== STEP 6: FEATURE-TRANSFER EXPERIMENT ==
=====
```

[INFO] Preparing flat features for MLP training...

[INFO] Initializing feedforward network...

[INFO] Training MLP on raw features...

Epoch 1/20	Train Loss: 0.5660	Acc: 0.7928
Epoch 2/20	Train Loss: 0.0612	Acc: 0.9992
Epoch 3/20	Train Loss: 0.0102	Acc: 0.9998
Epoch 4/20	Train Loss: 0.0038	Acc: 1.0000
Epoch 5/20	Train Loss: 0.0020	Acc: 1.0000
Epoch 6/20	Train Loss: 0.0012	Acc: 1.0000
Epoch 7/20	Train Loss: 0.0008	Acc: 1.0000
Epoch 8/20	Train Loss: 0.0006	Acc: 1.0000
Epoch 9/20	Train Loss: 0.0004	Acc: 1.0000

```

Epoch 10/20 | Train Loss: 0.0003 Acc: 1.0000
Epoch 11/20 | Train Loss: 0.0003 Acc: 1.0000
Epoch 12/20 | Train Loss: 0.0002 Acc: 1.0000
Epoch 13/20 | Train Loss: 0.0002 Acc: 1.0000
Epoch 14/20 | Train Loss: 0.0001 Acc: 1.0000
Epoch 15/20 | Train Loss: 0.0001 Acc: 1.0000
Epoch 16/20 | Train Loss: 0.0001 Acc: 1.0000
Epoch 17/20 | Train Loss: 0.0001 Acc: 1.0000
Epoch 18/20 | Train Loss: 0.0001 Acc: 1.0000
Epoch 19/20 | Train Loss: 0.0001 Acc: 1.0000
Epoch 20/20 | Train Loss: 0.0001 Acc: 1.0000

[INFO] Freezing feature extractor and generating embeddings...
[INFO] Train embeddings shape: (5281, 64)
[INFO] Test embeddings shape: (1321, 64)

[INFO] Training Logistic Regression on learned embeddings...
[RESULTS] Accuracy of Logistic Regression on learned embeddings: 1.0000

[INFO] Training Logistic Regression on raw features...
[RESULTS] Accuracy of Logistic Regression on raw features: 1.0000

=====
FEATURE-TRANSFER COMPARISON RESULTS
-----
Logistic Regression on raw features: 1.0000
Logistic Regression on learned features: 1.0000
=====
```

Definitely lessons learned about accuracy in this project.

1.000 = yeah, right.. something's off here.

Overfitting? Leakage? Hmm.

```

from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
import pandas as pd
import numpy as np
import torch
import torch.nn as nn
from torch.utils.data import DataLoader, TensorDataset
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt
import seaborn as sns

print("""
=====
==  STEP 6: FEATURE-TRANSFER EXPERIMENT (TWEAKED)      ==
=====
""")
```

```

# -----
# 1. Config
# -----
BATCH_SIZE = 32
EPOCHS = 20
LR = 1e-3
HIDDEN_SIZE = 64

# -----
# 2. Prepare features & target
# -----
flat_feature_cols = [
    'goal_difference', 'total_goals', 'expected_goal_diff',
    'home_avg_goal_diff_last5', 'away_avg_goal_diff_last5',
    'match_hour', 'match_day_of_week', 'match_month'
]

scaler = StandardScaler()
X_flat = scaler.fit_transform(cleaned_df[flat_feature_cols])

ohe_target = OneHotEncoder(sparse_output=False, handle_unknown='ignore')
y_flat = ohe_target.fit_transform(cleaned_df[['winner']])
target_classes = list(ohe_target.categories_[0])

# Time-based split option (comment out if not needed)
# split_idx = int(len(X_flat) * 0.8)
# X_train, X_test = X_flat[:split_idx], X_flat[split_idx:]
# y_train, y_test = y_flat[:split_idx], y_flat[split_idx:]

# Random split
X_train, X_test, y_train, y_test = train_test_split(
    X_flat, y_flat, test_size=0.2, random_state=42, shuffle=True
)

train_dataset = TensorDataset(torch.tensor(X_train, dtype=torch.float32),
                             torch.tensor(y_train, dtype=torch.float32))
test_dataset = TensorDataset(torch.tensor(X_test, dtype=torch.float32),
                            torch.tensor(y_test, dtype=torch.float32))

train_loader = DataLoader(train_dataset, batch_size=BATCH_SIZE, shuffle=True)
test_loader = DataLoader(test_dataset, batch_size=BATCH_SIZE)

# -----
# 3. Define MLP
# -----
class FeedforwardNet(nn.Module):
    def __init__(self, input_dim, hidden_dim, output_dim):
        super().__init__()
        self.feature_extractor = nn.Sequential(
            nn.Linear(input_dim, hidden_dim).

```

```

        ....)
        nn.ReLU(),
        nn.Linear(hidden_dim, hidden_dim),
        nn.ReLU()
    )
    self.classifier = nn.Linear(hidden_dim, output_dim)

def forward(self, x):
    feats = self.feature_extractor(x)
    out = self.classifier(feats)
    return out, feats

model = FeedforwardNet(X_train.shape[1], HIDDEN_SIZE, y_train.shape[1])
criterion = nn.CrossEntropyLoss()
optimizer = torch.optim.Adam(model.parameters(), lr=LR)

# -----
# 4. Train MLP with val check
# -----
print("\n[INFO] Training MLP on raw features...")
for epoch in range(EPOCHS):
    model.train()
    running_loss, correct, total = 0.0, 0, 0
    for X_batch, y_batch in train_loader:
        optimizer.zero_grad()
        outputs, _ = model(X_batch)
        loss = criterion(outputs, torch.argmax(y_batch, dim=1))
        loss.backward()
        optimizer.step()
        running_loss += loss.item() * X_batch.size(0)
        preds = torch.argmax(outputs, dim=1)
        correct += (preds == torch.argmax(y_batch, dim=1)).sum().item()
        total += y_batch.size(0)
    train_loss = running_loss / total
    train_acc = correct / total

    # Validation accuracy
    model.eval()
    val_correct, val_total = 0, 0
    with torch.no_grad():
        for X_batch, y_batch in test_loader:
            outputs, _ = model(X_batch)
            preds = torch.argmax(outputs, dim=1)
            val_correct += (preds == torch.argmax(y_batch, dim=1)).sum().item()
            val_total += y_batch.size(0)
    val_acc = val_correct / val_total

    print(f"Epoch {epoch+1}/{EPOCHS} | Train Loss: {train_loss:.4f} "
          f"Train Acc: {train_acc:.4f} | Val Acc: {val_acc:.4f}")

# -----

```

```

# 5. Freeze & extract embeddings
# -----
for param in model.feature_extractor.parameters():
    param.requires_grad = False

with torch.no_grad():
    train_embeddings = model.feature_extractor(torch.tensor(X_train, dtype=torch.float32))
    test_embeddings = model.feature_extractor(torch.tensor(X_test, dtype=torch.float32))

# -----
# 6. Logistic Regression on embeddings
# -----
log_reg = LogisticRegression(max_iter=1000)
log_reg.fit(train_embeddings, np.argmax(y_train, axis=1))
y_pred = log_reg.predict(test_embeddings)
transfer_acc = accuracy_score(np.argmax(y_test, axis=1), y_pred)

# -----
# 7. Logistic Regression on raw features
# -----
log_reg_baseline = LogisticRegression(max_iter=1000)
log_reg_baseline.fit(X_train, np.argmax(y_train, axis=1))
y_pred_base = log_reg_baseline.predict(X_test)
baseline_acc = accuracy_score(np.argmax(y_test, axis=1), y_pred_base)

# -----
# 8. Comparison & diagnostics
# -----
print("\n=====")
print(" FEATURE-TRANSFER COMPARISON RESULTS")
print("-----")
print(f" Logistic Regression on raw features: {baseline_acc:.4f}")
print(f" Logistic Regression on learned features: {transfer_acc:.4f}")
print("=====")

print("\n[CLASSIFICATION REPORT - Raw Features]")
print(classification_report(np.argmax(y_test, axis=1), y_pred_base, target_names=target_classes))

print("\n[CLASSIFICATION REPORT - Learned Features]")
print(classification_report(np.argmax(y_test, axis=1), y_pred, target_names=target_classes))

# Confusion matrix plot for learned features
cm = confusion_matrix(np.argmax(y_test, axis=1), y_pred)
plt.figure(figsize=(5,4))
sns.heatmap(cm, annot=True, fmt="d", cmap="Blues",
            xticklabels=target_classes, yticklabels=target_classes)
plt.xlabel("Predicted")
plt.ylabel("True")
plt.title("Confusion Matrix - Logistic Regression on Learned Features")
plt.show()

```


ARGH! Still not right. It's gotta be that 'goal difference' feature again -- it's fine for overall stats, but doesn't help in predictive analysis. One more run with pre-match data and rolling averages should do the trick.

```
#####
# ### STEP 6: FEATURE-TRANSFER EXPERIMENT – LEAKAGE-SAFE VERSION ###
#####

print("""
=====
== STEP 6: FEATURE-TRANSFER EXPERIMENT (LEAKAGE-SAFE) ==
=====

""")

import pandas as pd
import numpy as np
import torch
import torch.nn as nn
from torch.utils.data import DataLoader, TensorDataset
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report, confusion_ma
```

```

from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt
import seaborn as sns

# =====
# === 1. CONFIGURATION =====
# =====
BATCH_SIZE = 32
EPOCHS = 20
LR = 1e-3
HIDDEN_SIZE = 64

# =====
# === 2. PREPARE SAFE FEATURES & TARGET =====
# =====
print("\n[INFO] Preparing leakage-safe features for MLP training...")

# Only pre-match features – no current match final stats
safe_feature_cols = [
    'expected_goal_diff', # if pre-match estimate
    'home_avg_goal_diff_last5',
    'away_avg_goal_diff_last5',
    'match_hour',
    'match_day_of_week',
    'match_month'
]

# Scale numeric features
scaler = StandardScaler()
X_flat = scaler.fit_transform(cleaned_df[safe_feature_cols])

# One-hot encode target
ohe_target = OneHotEncoder(sparse_output=False, handle_unknown='ignore')
y_flat = ohe_target.fit_transform(cleaned_df[['winner']])
target_classes = list(ohe_target.categories_[0])

# Train/test split (random; switch to time-based if needed)
X_train, X_test, y_train, y_test = train_test_split(
    X_flat, y_flat, test_size=0.2, random_state=42, shuffle=True
)

train_dataset = TensorDataset(torch.tensor(X_train, dtype=torch.float32),
                             torch.tensor(y_train, dtype=torch.float32))
test_dataset = TensorDataset(torch.tensor(X_test, dtype=torch.float32),
                           torch.tensor(y_test, dtype=torch.float32))

train_loader = DataLoader(train_dataset, batch_size=BATCH_SIZE, shuffle=True)
test_loader = DataLoader(test_dataset, batch_size=BATCH_SIZE)

# =====

```

```

# === 3. DEFINE FEEDFORWARD NETWORK =====
# =====
class FeedforwardNet(nn.Module):
    def __init__(self, input_dim, hidden_dim, output_dim):
        super().__init__()
        self.feature_extractor = nn.Sequential(
            nn.Linear(input_dim, hidden_dim),
            nn.ReLU(),
            nn.Linear(hidden_dim, hidden_dim),
            nn.ReLU()
        )
        self.classifier = nn.Linear(hidden_dim, output_dim)

    def forward(self, x):
        feats = self.feature_extractor(x)
        out = self.classifier(feats)
        return out, feats

model = FeedforwardNet(X_train.shape[1], HIDDEN_SIZE, y_train.shape[1])
criterion = nn.CrossEntropyLoss()
optimizer = torch.optim.Adam(model.parameters(), lr=LR)

# =====
# === 4. TRAIN MLP ON SAFE FEATURES =====
# =====
print("\n[INFO] Training MLP on leakage-safe features...")
for epoch in range(EPOCHS):
    model.train()
    running_loss, correct, total = 0.0, 0, 0
    for X_batch, y_batch in train_loader:
        optimizer.zero_grad()
        outputs, _ = model(X_batch)
        loss = criterion(outputs, torch.argmax(y_batch, dim=1))
        loss.backward()
        optimizer.step()
        running_loss += loss.item() * X_batch.size(0)
        preds = torch.argmax(outputs, dim=1)
        correct += (preds == torch.argmax(y_batch, dim=1)).sum().item()
        total += y_batch.size(0)
    train_loss = running_loss / total
    train_acc = correct / total

    # Validation accuracy
    model.eval()
    val_correct, val_total = 0, 0
    with torch.no_grad():
        for X_batch, y_batch in test_loader:
            outputs, _ = model(X_batch)
            preds = torch.argmax(outputs, dim=1)
            val_correct += (preds == torch.argmax(y_batch, dim=1)).sum().item()
            val_total += y_batch.size(0)

```

```

    val_total += y_val.size(0)
    val_acc = val_correct / val_total

    print(f"Epoch {epoch+1}/{EPOCHS} | Train Loss: {train_loss:.4f} "
          f"Train Acc: {train_acc:.4f} | Val Acc: {val_acc:.4f}")

# =====
# === 5. FREEZE & EXTRACT EMBEDDINGS =====
# =====
for param in model.feature_extractor.parameters():
    param.requires_grad = False

with torch.no_grad():
    train_embeddings = model.feature_extractor(torch.tensor(X_train, dtype=torch.float32))
    test_embeddings = model.feature_extractor(torch.tensor(X_test, dtype=torch.float32))

# =====
# === 6. LOGISTIC REGRESSION COMPARISONS =====
# =====
log_reg = LogisticRegression(max_iter=1000)
log_reg.fit(train_embeddings, np.argmax(y_train, axis=1))
y_pred = log_reg.predict(test_embeddings)
transfer_acc = accuracy_score(np.argmax(y_test, axis=1), y_pred)

log_reg_baseline = LogisticRegression(max_iter=1000)
log_reg_baseline.fit(X_train, np.argmax(y_train, axis=1))
y_pred_base = log_reg_baseline.predict(X_test)
baseline_acc = accuracy_score(np.argmax(y_test, axis=1), y_pred_base)

# =====
# === 7. OUTPUT & DIAGNOSTICS =====
# =====
print("\n=====")
print(" FEATURE-TRANSFER COMPARISON RESULTS (LEAKAGE-SAFE)")
print("-----")
print(f" Logistic Regression on raw safe features: {baseline_acc:.4f}")
print(f" Logistic Regression on learned safe features: {transfer_acc:.4f}")
print("=====")

print("\n[CLASSIFICATION REPORT - Raw Safe Features]")
print(classification_report(np.argmax(y_test, axis=1), y_pred_base, target_names=target_classes))

print("\n[CLASSIFICATION REPORT - Learned Safe Features]")
print(classification_report(np.argmax(y_test, axis=1), y_pred, target_names=target_classes))

# Confusion matrix plot for learned features
cm = confusion_matrix(np.argmax(y_test, axis=1), y_pred)
plt.figure(figsize=(5,4))
sns.heatmap(cm, annot=True, fmt="d", cmap="Blues",
            xticklabels=target_classes, yticklabels=target_classes)
plt.xlabel("Predicted")

```

```
plt.ylabel("True")
plt.title("Confusion Matrix – Logistic Regression on Learned Safe Features")
plt.show()
```

▼ EVAL ON RNN v MLP/LR

Great, neither model is now "cheating" with leakage or other issues that skew the data. Accuracy is still more like a coin flip.

It's fascinating how draws are avoided, to a degree.. almost like bettors and bookmakers don't like to deal with ties. It's probably because draws can result in a lot of other factors that, initially, are hard to see *before* the match starts.

But let's run a calculated comparison and see..

```
# =====
# === MODEL COMPARISON HEATMAP + COMMENTS (LEAKAGE-SAFE) ====
# =====
```

```

import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Replace these with the actual values from your runs
# (Here I'm using the leakage-safe example numbers you posted earlier)
df_compare = pd.DataFrame([
    {"Model": "RNN (LSTM, safe features)", "Accuracy": 0.408333, "Macro-F1": 0.21},
    {"Model": "MLP (safe features)", "Accuracy": 0.481500, "Macro-F1": 0.41},
    {"Model": "Kernel Regressor (LR)", "Accuracy": 0.466300, "Macro-F1": 0.31}
])

# --- Plot heatmap ---
plt.figure(figsize=(8, 4))
sns.heatmap(df_compare.set_index("Model"), annot=True, fmt=".3f", cmap="Blues")
plt.title("Model Performance Comparison (Leakage-Safe)", fontsize=14, pad=12)
plt.yticks(rotation=0)
plt.show()

# --- Print commentary ---
best_acc_model = df_compare.loc[df_compare['Accuracy'].idxmax(), 'Model']
best_f1_model = df_compare.loc[df_compare['Macro-F1'].idxmax(), 'Model']
fastest_model = df_compare.loc[df_compare['Train Time (s)'].idxmin(), 'Model']

print("\n[SUMMARY] Model Performance Comparison (Leakage-Safe):")
print(f"- Best Overall Accuracy: {best_acc_model} ({df_compare['Accuracy'].max()}")
print(f"- Best Macro-F1 (class balance): {best_f1_model} ({df_compare['Macro-F1']}")
print(f"- Fastest to Train: {fastest_model} ({df_compare['Train Time (s)'].min()")

print("\n[COMMENTS]")
print("• MLP (safe features) leads in both accuracy and macro-F1, showing better")
print("• Kernel Regressor (LR) is nearly as accurate as the MLP, but trains almost")
print("• RNN (LSTM) lags in both accuracy and macro-F1, suggesting it's not leveraging")
print("• All models still struggle with draws – macro-F1 scores show imbalance in the")

```



✓ YES THF Central Park Squirrel Census!

DATASET TWO!

Why?

Size: 3,023 sightings (JUST meets the project's minimum)

Features: Location coordinates, age, fur color, elevation, activities, communications, interactions

Read more here: https://data.cityofnewyork.us/Environment/2018-Central-Park-Squirrel-Census-Squirrel-Data/vfnx-vebw/about_data

Modeling Ideas:

Classification: Predict squirrel behavior type (e.g. foraging, climbing) based on physical traits and location

Regression: Estimate likelihood of human interaction or activity level

RNN: If you group sightings by time or location, you could model behavioral sequences (e.g. "What does a squirrel do next?")

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer

print(" --- Starting Squirrel Data ---")

# --- Step 1: Data Gathering & Initial Inspection ---
# Direct URL to the raw CSV file
data_url = 'https://data.cityofnewyork.us/api/views/vfnx-vebw/rows.csv?accessType=CSV'
squirrel_df = pd.read_csv(data_url)

# Display basic info and first few rows to get a feel for the data
print("\nInitial Data Info:")
print(squirrel_df.info())
print("\nFirst 5 Rows:")
print(squirrel_df.head().to_string())
```

--- Starting Squirrel Data ---

Initial Data Info:

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 3023 entries, 0 to 3022

Data columns (total 31 columns):

#	Column	Non-Null Count	Dtype
0	X	3023 non-null	float64
1	Y	3023 non-null	float64
2	Unique Squirrel ID	3023 non-null	object
3	Hectare	3023 non-null	object
4	Shift	3023 non-null	object
5	Date	3023 non-null	int64
6	Hectare Squirrel Number	3023 non-null	int64
7	Age	2902 non-null	object
8	Primary Fur Color	2968 non-null	object
9	Highlight Fur Color	1937 non-null	object
10	Combination of Primary and Highlight Color	3023 non-null	object
11	Color notes	182 non-null	object
12	Location	2959 non-null	object
13	Above Ground Sighter Measurement	2909 non-null	object
14	Specific Location	476 non-null	object
15	Running	3023 non-null	bool
16	Chasing	3023 non-null	bool
17	Climbing	3023 non-null	bool
18	Eating	3023 non-null	bool
19	Foraging	3023 non-null	bool
20	Other Activities	437 non-null	object
21	Kuks	3023 non-null	bool
22	Quaas	3023 non-null	bool
23	Moans	3023 non-null	bool
24	Tail flags	3023 non-null	bool
25	Tail twitches	3023 non-null	bool
26	Approaches	3023 non-null	bool
27	Indifferent	3023 non-null	bool
28	Runs from	3023 non-null	bool
29	Other Interactions	240 non-null	object
30	Lat/Long	3023 non-null	object

dtypes: bool(13), float64(2), int64(2), object(14)

memory usage: 463.6+ KB

None

First 5 Rows:

	X	Y	Unique Squirrel ID	Hectare	Shift	Date	Hectare Squi
0	-73.956134	40.794082	37F-PM-1014-03	37F	PM	10142018	
1	-73.968857	40.783783	21B-AM-1019-04	21B	AM	10192018	
2	-73.974281	40.775534	11B-PM-1014-08	11B	PM	10142018	
3	-73.959641	40.790313	32E-PM-1017-14	32E	PM	10172018	
4	-73.970268	40.776213	13E-AM-1017-05	13E	AM	10172018	

```
# STEP 1 - part 2 - FOR SECOND DATASET
# Check for null values in each column
print("\nMissing values per column:")
```

```
print(squirrel_df.isnull().sum())
```

```
Missing values per column:  
X 0  
Y 0  
Unique Squirrel ID 0  
Hectare 0  
Shift 0  
Date 0  
Hectare Squirrel Number 0  
Age 121  
Primary Fur Color 55  
Highlight Fur Color 1086  
Combination of Primary and Highlight Color 0  
Color notes 2841  
Location 64  
Above Ground Sighter Measurement 114  
Specific Location 2547  
Running 0  
Chasing 0  
Climbing 0  
Eating 0  
Foraging 0  
Other Activities 2586  
Kuks 0  
Quaas 0  
Moans 0  
Tail flags 0  
Tail twitches 0  
Approaches 0  
Indifferent 0  
Runs from 0  
Other Interactions 2783  
Lat/Long 0  
dtype: int64
```

✓ DECISION TIME!

Based on the provided missing value counts, here's a recommended strategy for which columns to impute or drop, with a justification for each choice.

Columns to Drop

It's best to drop columns with a very high percentage of missing values, as they are unlikely to provide meaningful information and could introduce noise into your analysis.

Color notes (2,841 missing)

Specific Location (2,547 missing)

Other Activities (2,586 missing)

Other Interactions (2,783 missing)

These columns are largely empty. For example, the high number of missing values for Other Interactions likely means that no "other" interactions were observed for most squirrels, making the column a poor candidate for imputation and of little use for the project.

Columns to Impute

For columns with a manageable number of missing values, imputation is the best strategy to retain valuable data. The choice of imputation method (e.g., mode for categorical, median for numerical) depends on the column type.

Age (121 missing): This is a key categorical feature. You can impute these missing

Primary Fur Color (55 missing): Another important categorical feature. Imputing wi

Highlight Fur Color (1,086 missing): While it has a high number of missing values,

Location (64 missing): This is a categorical field. You can impute the missing value

Above Ground Sighter Measurement (114 missing): This is a numerical feature. You can

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer

print("--- Starting Squirrel Data Missing Values ---")

# --- Step 2: Data Cleaning and Imputation ---
# Drop columns with too many missing values, as previously planned
columns_to_drop = [
    "Color notes",
    "Specific Location",
```

```

        "Other Activities",
        "Other Interactions"
    ]
squirrel_df = squirrel_df.drop(columns=columns_to_drop, errors='ignore')

# Impute missing values based on the notebook's strategy
# Fill 'Highlight Fur Color' with a new category 'No Highlight'
squirrel_df['Highlight Fur Color'] = squirrel_df['Highlight Fur Color'].fillna("No Highlight")

# Use mode imputation for other categorical features
for col in ["Age", "Primary Fur Color", "Location"]:
    mode_val = squirrel_df[col].mode()
    if not mode_val.empty:
        squirrel_df[col].fillna(mode_val[0], inplace=True)
    else:
        print(f"Warning: Could not find mode for '{col}' column.")

# Convert 'Above Ground Sighter Measurement' to numeric and impute with median
squirrel_df['Above Ground Sighter Measurement'] = pd.to_numeric(squirrel_df['Above Ground Sighter Measurement'])
median_val = squirrel_df['Above Ground Sighter Measurement'].median()
squirrel_df['Above Ground Sighter Measurement'].fillna(median_val, inplace=True)

# Final check of the cleaned data
print("\nCleaned Data Info:")
print(squirrel_df.info())

```

--- Starting Squirrel Data Missing Values ---

Cleaned Data Info:

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 3023 entries, 0 to 3022

Data columns (total 27 columns):

#	Column	Non-Null Count	Dtype
0	X	3023 non-null	float64
1	Y	3023 non-null	float64
2	Unique Squirrel ID	3023 non-null	object
3	Hectare	3023 non-null	object
4	Shift	3023 non-null	object
5	Date	3023 non-null	int64
6	Hectare Squirrel Number	3023 non-null	int64
7	Age	3023 non-null	object
8	Primary Fur Color	3023 non-null	object
9	Highlight Fur Color	3023 non-null	object
10	Combination of Primary and Highlight Color	3023 non-null	object
11	Location	3023 non-null	object
12	Above Ground Sighter Measurement	3023 non-null	float64
13	Running	3023 non-null	bool
14	Chasing	3023 non-null	bool
15	Climbing	3023 non-null	bool
16	Eating	3023 non-null	bool
17	Foraging	3023 non-null	bool
18	Kuks	3023 non-null	bool

```
-- 19 Quaas 3023 non-null bool  
-- 20 Moans 3023 non-null bool  
-- 21 Tail flags 3023 non-null bool  
-- 22 Tail twitches 3023 non-null bool  
-- 23 Approaches 3023 non-null bool  
-- 24 Indifferent 3023 non-null bool  
-- 25 Runs from 3023 non-null bool  
-- 26 Lat/Long 3023 non-null object  
dtypes: bool(13), float64(3), int64(2), object(9)  
memory usage: 369.1+ KB  
None  
/tmp/ipython-input-1380712073.py:33: FutureWarning: A value is trying to be set  
The behavior will change in pandas 3.0. This inplace method will never work bca
```

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.met

```
squirrel_df[col].fillna(mode_val[0], inplace=True)  
/tmp/ipython-input-1380712073.py:40: FutureWarning: A value is trying to be set  
The behavior will change in pandas 3.0. This inplace method will never work bca
```

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.met

```
squirrel_df['Above Ground Sighter Measurement'].fillna(median_val, inplace=True)
```

We loaded the NYC squirrel dataset and cleaned it by removing columns that had too many blanks, like "Color notes" and "Other Interactions." For the other missing values, we filled them in so the dataset stayed complete. For example, we used the most common value for things like "Age," "Primary Fur Color," and "Location," while for "Highlight Fur Color" we added "No Highlight" instead of leaving it empty. For the numeric "Above Ground Sighter Measurement," we used the median to fill gaps. After this, the dataset had no missing values and was ready to use.

```
import pandas as pd  
import numpy as np  
import matplotlib.pyplot as plt  
import seaborn as sns  
from sklearn.model_selection import train_test_split  
from sklearn.linear_model import LogisticRegression  
from sklearn.metrics import accuracy_score, classification_report  
from sklearn.preprocessing import OneHotEncoder, StandardScaler  
from sklearn.compose import ColumnTransformer  
from sklearn.pipeline import Pipeline  
from sklearn.impute import SimpleImputer  
  
print("--- Starting Squirrel Data EDA ---")  
  
# --- Step 3: Exploratory Data Analysis (EDA) and Visualizations ---
```

```

print("\n--- Performing EDA and Generating Visualizations ---")

# Plot distribution of the target variable 'Running'
plt.figure(figsize=(6, 4))
sns.countplot(x='Running', data=squirrel_df)
plt.title('Distribution of Running vs. Not Running Squirrels')
plt.xlabel('Is Running?')
plt.ylabel('Count')
plt.show()
print("Interpretation: This shows how balanced or imbalanced our target variable is. The count for 'Not Running' is approximately 1000, while 'Running' is approximately 200, indicating an imbalance where most squirrels are not running.")

# Plot distribution of the main categorical feature: 'Primary Fur Color'
plt.figure(figsize=(10, 6))
sns.countplot(x='Primary Fur Color', data=squirrel_df, order=squirrel_df['Primary Fur Color'].unique())
plt.title('Distribution of Primary Fur Color')
plt.xlabel('Primary Fur Color')
plt.ylabel('Count')
plt.show()
print("Interpretation: This helps us understand the most common fur colors in our dataset. The most frequent color is Gray, followed by Black, and then Brown and Cinnamon.")


# Plot the relationship between a numerical feature and the target variable
plt.figure(figsize=(8, 6))
sns.boxplot(x='Running', y='Above Ground Sighter Measurement', data=squirrel_df)
plt.title('Above Ground Sighter Measurement by Running Behavior')
plt.xlabel('Is Running?')
plt.ylabel('Above Ground Measurement (ft)')
plt.show()
print("Interpretation: We can see if the height of the squirrel is related to its running behavior. The median height for running squirrels is slightly higher than for non-running ones, though there is significant overlap.")


# Scatter plot of location
plt.figure(figsize=(10, 8))
plt.scatter(squirrel_df['X'], squirrel_df['Y'], alpha=0.5, c=squirrel_df['Running'])
plt.title('Squirrel Sightings in Central Park')
plt.xlabel('Longitude (X)')
plt.ylabel('Latitude (Y)')
plt.legend(handles=[
    plt.Line2D([0], [0], marker='o', color='w', label='Not Running', markerfacecolor='white'),
    plt.Line2D([0], [0], marker='o', color='w', label='Running', markerfacecolor='white')
])
plt.show()
print("Interpretation: This gives a spatial overview of where sightings occurred. Red dots represent non-running squirrels, and blue dots represent running squirrels. There appears to be a higher density of sightings in the central part of the park.")


# --- Step 4: Outlier Detection ---
col = "Above Ground Sighter Measurement"
Q1 = squirrel_df[col].quantile(0.25)
Q3 = squirrel_df[col].quantile(0.75)
IQR = Q3 - Q1
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR
outliers = squirrel_df[(squirrel_df[col] < lower_bound) | (squirrel_df[col] > upper_bound)]
print(f"Number of outliers: {len(outliers)}")

```

```
print(f"\nnumber of outliers in '{col}': {len(outliers)}")\n\n# --- Step 5: Feature Encoding and Scaling (Pre-modeling step) ---\ntarget_variable = 'Running'\nX_squirrel = squirrel_df.drop(columns=[target_variable])\ny_squirrel = squirrel_df[target_variable].astype(int)\n\nnumerical_features = X_squirrel.select_dtypes(include=np.number).columns.tolist()\ncategorical_features = X_squirrel.select_dtypes(include=['object', 'category']).columns.tolist()\n\n# Define the preprocessing pipeline\nnumerical_pipeline = Pipeline(steps=[\n    ('imputer', SimpleImputer(strategy='median')),\n    ('scaler', StandardScaler())\n])\n\npreprocessor = ColumnTransformer(\n    transformers=[\n        ('num', numerical_pipeline, numerical_features),\n        ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features)\n    ],\n    remainder='passthrough'\n)\n\n# Apply preprocessing to all data (will be split later for modeling)\nX_processed_squirrel = preprocessor.fit_transform(X_squirrel)\nprint("\nFeatures have been successfully encoded and scaled, ready for modeling.\nFinal feature matrix shape: {X_processed_squirrel.shape}\")\n\nprint("\n--- EDA and Preprocessing Complete! ---")
```



```

print("\n--- Feature Engineering (AKA stuff that makes sense for squirrels) ---")

# ok so... first idea: sum up all the "yes/no" behavior flags into one number
# more active squirrels are more likely to run, right?
behavioral_flags = [
    'Chasing', 'Climbing', 'Eating', 'Foraging', 'Kuks', 'Quaas', 'Moans',
    'Tail flags', 'Tail twitches', 'Approaches', 'Indifferent', 'Runs from'
]
squirrel_df['activity_intensity'] = squirrel_df[behavioral_flags].astype(int).sum()
print("activity_intensity added. sample:", squirrel_df['activity_intensity'][:5])

# second: how far from the "middle" of Central Park is this squirrel?
# Do squirrels near the edges act different? Probably if there's less trees.. maybe
x0 = squirrel_df['X'].median()
y0 = squirrel_df['Y'].median()
squirrel_df['dist_from_center'] = np.sqrt((squirrel_df['X'] - x0)**2 + (squirrel_df['Y'] - y0)**2)
print("dist_from_center added. median center was:", (x0, y0))

# third: put heights into quartiles
# maybe "higher squirrels" aren't running.. because a fall would hurt from 100 feet
try:
    squirrel_df['height_class'] = pd.qcut(
        squirrel_df['Above Ground Sighter Measurement'],
        q=4, labels=False, duplicates='drop')
except Exception as e:
    print("uh oh, height_class binning failed:", e)

# quick peek at what we just did
print("\nHEAD of engineered features:")
print(squirrel_df[['activity_intensity', 'dist_from_center', 'height_class']].head)

# --- VISUALS for domain---
import matplotlib.pyplot as plt
import seaborn as sns

# Activity intensity distribution
plt.figure(figsize=(6,4))
sns.countplot(x='activity_intensity', data=squirrel_df)
plt.title('Distribution of Activity Intensity')

```

```
plt.title('Distribution of Activity Intensity')
plt.xlabel('Number of Active Behaviors')
plt.ylabel('Count')
plt.show()

# Dist from center histogram
plt.figure(figsize=(6,4))
sns.histplot(squirrel_df['dist_from_center'], bins=30, kde=True)
plt.title('Distance from Park Center')
plt.xlabel('Distance (map units)')
plt.ylabel('Count')
plt.show()

# Height class vs Running
plt.figure(figsize=(6,4))
sns.countplot(x='height_class', hue='Running', data=squirrel_df)
plt.title('Height Class vs Running Behavior')
plt.xlabel('Height Class (quartiles)')
plt.ylabel('Count')
plt.legend(title='Running')
plt.show()

# --- Quick stats ---
print("\nFeature summary stats:")
print(squirrel_df[['activity_intensity','dist_from_center','height_class']].desc()
```



```

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

# --- NOTE: This code assumes squirrel_df has been loaded and cleaned as per previous
# If not, you must run the data loading and cleaning code first.

col = "Above Ground Sighter Measurement"

# --- Visualization 1: Box Plot to show quartiles and outliers ---
plt.figure(figsize=(8, 6))
sns.boxplot(y=squirrel_df[col])
plt.title(f'Box Plot of {col}')
plt.ylabel(col)
plt.show()

print("Interpretation: The box represents the middle 50% of the data (IQR), the whiskers represent the full range of the data, and the outliers are points outside the whisker boundaries.")

# --- Visualization 2: Histogram with Outlier Bounds ---
plt.figure(figsize=(10, 6))
sns.histplot(squirrel_df[col], bins=30, kde=True)

# Calculate outlier bounds (re-calculated for clarity)
Q1 = squirrel_df[col].quantile(0.25)
Q3 = squirrel_df[col].quantile(0.75)
IQR = Q3 - Q1
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

# Add vertical lines to show the outlier boundaries
plt.axvline(x=lower_bound, color='r', linestyle='--', label='Lower Bound')
plt.axvline(x=upper_bound, color='r', linestyle='--', label='Upper Bound')

plt.title(f'Distribution of {col} with Outlier Bounds')
plt.xlabel(col)
plt.ylabel('Count')
plt.legend()
plt.show()

```

```
print(f"Interpretation: The values outside the red dashed lines are considered ou
```

```

import numpy as np

# Select the numeric column
col = "Above Ground Sighter Measurement"

# Calculate Q1 (25th percentile) and Q3 (75th percentile)
Q1 = squirrel_df[col].quantile(0.25)
Q3 = squirrel_df[col].quantile(0.75)
IQR = Q3 - Q1

# Define outlier boundaries
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

# Find outliers
outliers = squirrel_df[(squirrel_df[col] < lower_bound) | (squirrel_df[col] > upper_bound)]

print(f"Number of outliers in {col}: {len(outliers)}")
print(f"Lower bound: {lower_bound}, Upper bound: {upper_bound}")
print("Example outliers:")
print(outliers.head())

```

Number of outliers in Above Ground Sighter Measurement: 677

Lower bound: 10.0, Upper bound: 10.0

Example outliers:

	X	Y	Unique Squirrel ID	Hectare	Shift	Date	\
9	-73.972250	40.774288	11D-AM-1010-03	11D	AM	10102018	
15	-73.970378	40.778753	16C-PM-1018-03	16C	PM	10182018	
19	-73.967113	40.778486	17F-AM-1007-07	17F	AM	10072018	
33	-73.971615	40.781391	18A-PM-1018-01	18A	PM	10182018	
38	-73.957465	40.789251	31H-PM-1008-02	31H	PM	10082018	

	Hectare	Squirrel Number	Age	Primary Fur Color	Highlight Fur Color	\
9		3	Adult	Gray	Cinnamon	
15		3	Adult	Gray	Cinnamon, White	
19		7	Adult	Gray	White	
33		1	Adult	Gray	No Highlight	
38		2	Juvenile	Gray	Cinnamon	

... Moans Tail flags Tail twitches Approaches Indifferent Runs from \

```

9   ...  means  null_flags  null_swiftlets  approachable  inactivity_rate  runs  vissim \
15  ...  False    False        False        False        True      False
19  ...  False    False        False        False        True      False
33  ...  False    False        False        False        False     False
38  ...  False    False        True         False        True      False

                                Lat/Long  activity_intensity \
9     POINT (-73.9722500196844 40.7742879599026)           2
15    POINT (-73.9703781726172 40.7787526130321)           3
19    POINT (-73.9671130680114 40.7784859700171)           2
33    POINT (-73.9716147061553 40.7813911036179)           1
38    POINT (-73.9574648097543 40.78925084286221)           3

            dist_from_center  height_class
9             0.005330          1
15            0.001878          0
19            0.001515          1
33            0.004419          1
38            0.015708          0

[5 rows x 30 columns]

```

When we looked for unusual values in the “Above Ground Sighter Measurement” column, we used the IQR method, which basically checks how far numbers stray from the typical range. This helps us spot squirrels that were recorded at very unusual heights compared to the rest. Outliers like this might just be rare but real cases (like a squirrel seen very high up in a tree), so let's check how many there are and how extreme they look first.

```

import numpy as np
import pandas as pd

# NOTE: Assumes squirrel_df has been loaded and cleaned in a previous cell

# Calculate IQR for the Above Ground Sighter Measurement
Q1 = squirrel_df["Above Ground Sighter Measurement"].quantile(0.25)
Q3 = squirrel_df["Above Ground Sighter Measurement"].quantile(0.75)
IQR = Q3 - Q1

# Define bounds
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

print("--- Outlier Detection for 'Above Ground Sighter Measurement' ---")
print(f"Lower bound: {lower_bound:.2f}, Upper bound: {upper_bound:.2f}")

# The following code is for identifying, NOT removing, the outliers.
outliers = squirrel_df[
    (squirrel_df["Above Ground Sighter Measurement"] < lower_bound) |
    (squirrel_df["Above Ground Sighter Measurement"] > upper_bound)
]

```

```

]

print(f"Number of outliers detected: {len(outliers)}")
print(f"Total dataset size: {squirrel_df.shape[0]}")
print("\nDecision: Outliers will be retained for modeling as they represent")
print("legitimate, though rare, behaviors in a wild animal population.")

--- Outlier Detection for 'Above Ground Sighter Measurement' ---
Lower bound: 10.00, Upper bound: 10.00
Number of outliers detected: 677
Total dataset size: 3023

Decision: Outliers will be retained for modeling as they represent
legitimate, though rare, behaviors in a wild animal population.

```

We checked for extreme values in the “Above Ground Sighter Measurement” column using the interquartile range (IQR) method, which flags numbers that are way higher or lower than the rest.

▼ **Removing outliers here would be absurd. Squirrels climb a lot.**

The outliers may be higher than most trees in Central Park, but what about statutes? Structures? Or a large stage for a concert? Those are significantly larger than a tree, and could easily reach over 100 feet - so these outliers may be quite legitimate.

```

import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer

# NOTE: Assumes the initial data loading and cleaning has been run, creating `sq

print("--- Starting Squirrel Data Preprocessing for Modeling ---")

# Define the target variable and features
target_variable = 'Running'
X_squirrel = squirrel_df.drop(columns=[target_variable])
y_squirrel = squirrel_df[target_variable].astype(int)

# Identify feature types

```

```

# Identify feature types
numerical_features = X_squirrel.select_dtypes(include=np.number).columns.tolist()
categorical_features = X_squirrel.select_dtypes(include=['object', 'category']).columns.tolist()

# Define the preprocessing pipeline
numerical_pipeline = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='median')),
    ('scaler', StandardScaler())
])

preprocessor = ColumnTransformer(
    transformers=[
        ('num', numerical_pipeline, numerical_features),
        ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features)
    ],
    remainder='passthrough'
)

# Apply preprocessing to all data
X_processed_squirrel = preprocessor.fit_transform(X_squirrel)

print("\nFeatures have been successfully encoded and scaled, ready for modeling.")
print(f"Final feature matrix shape: {X_processed_squirrel.shape}")

--- Starting Squirrel Data Preprocessing for Modeling ---

Features have been successfully encoded and scaled, ready for modeling.
Final feature matrix shape: (3023, 6443)

```

✓ Um.. no!

That just seems like a bizarre shape. And that was verified with help from Gemini. The One Hot Encoding is treating something like a unique squirrel ID number or some other feature like it's a category, not an identifier.

So it's counting that as well, which could really screw things up. So we're going to exclude those items that can really alter things poorly.

```

import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer

print("--- Starting Corrected Squirrel Data Preprocessing ---")

```

```

# NOTE: Assumes the initial data loading and cleaning has been run, creating `sq

# Define the target variable and features
target_variable = 'Running'
X_squirrel = squirrel_df.drop(columns=[target_variable])
y_squirrel = squirrel_df[target_variable].astype(int)

# Identify feature types, EXCLUDING the identifier columns that caused the issue
numerical_features = ["X", "Y", "Above Ground Sighter Measurement"]
categorical_features = ["Age", "Primary Fur Color", "Highlight Fur Color", "Loca

# The boolean columns should be treated as numeric (0 or 1) and will not be one-l
boolean_cols = [
    'Chasing', 'Climbing', 'Eating', 'Foraging', 'Kuks', 'Quaas',
    'Moans', 'Tail flags', 'Tail twitches', 'Approaches', 'Indifferent',
    'Runs from'
]
numerical_features.extend(boolean_cols)

# Define the preprocessing pipeline
numerical_pipeline = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='median')),
    ('scaler', StandardScaler())
])

preprocessor = ColumnTransformer(
    transformers=[
        ('num', numerical_pipeline, numerical_features),
        ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features)
    ],
    remainder='drop' # Explicitly drop columns not in our lists
)

# Apply preprocessing to all data
X_processed_squirrel = preprocessor.fit_transform(X_squirrel)

print("\nFeatures have been successfully encoded and scaled, ready for modeling.")
print(f"Final feature matrix shape: {X_processed_squirrel.shape}")

```

--- Starting Corrected Squirrel Data Preprocessing ---

Features have been successfully encoded and scaled, ready for modeling.
 Final feature matrix shape: (3023, 36)

✓ MUCH BETTER!

That a shape we can work with.

When we ran the incorrect preprocessing that created over 6,000 features, it resulted in

When we ran the incorrect preprocessing that created over 6,000 features, it resulted in the following:

Old Classical Model Comparison for Squirrel Data

1. Random Forest 0.825532
2. Gradient Boosting 0.821277
3. XGBoost 0.817021
4. SVM (RBF Kernel) 0.797872
5. Logistic Regression 0.787234
6. K-Nearest Neighbors 0.765957

"Running" is fine as a target.. but kinda boring. :)

Let's try a couple more..

```
print("\n--- Setting up multiple classification targets ---")

# Define the targets we want to try
classification_targets = ['Running', 'Climbing', 'Chasing']

# Store them in a dict for easy looping later
target_data = {}
for target in classification_targets:
    y = squirrel_df[target].astype(int)
    target_data[target] = train_test_split(
        X_processed_squirrel, y,
        test_size=0.2, random_state=42, stratify=y
    )
    print(f"Target '{target}' ready. Class balance:\n{y.value_counts(normalize=True)}")
```

```
--- Setting up multiple classification targets ---
Target 'Running' ready. Class balance:
Running
0    0.758518
1    0.241482
Name: proportion, dtype: float64

Target 'Climbing' ready. Class balance:
Climbing
0    0.782335
1    0.217665
Name: proportion, dtype: float64

Target 'Chasing' ready. Class balance:
Chasing
0    0.907708
1    0.092292
```

```
+-----  
Name: proportion, dtype: float64
```

```
from sklearn.linear_model import LogisticRegression  
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier  
from sklearn.svm import SVC  
from sklearn.neighbors import KNeighborsClassifier  
from sklearn.metrics import accuracy_score  
import pandas as pd  
  
print("\n--- Training Classical Models for Multiple Targets ---")  
  
# Define the models once  
models = [  
    ("Logistic Regression", LogisticRegression(max_iter=1000, random_state=42)),  
    ("Random Forest", RandomForestClassifier(random_state=42)),  
    ("SVM (RBF Kernel)", SVC(kernel='rbf', random_state=42)),  
    ("Gradient Boosting", GradientBoostingClassifier(random_state=42)),  
    ("K-Nearest Neighbors", KNeighborsClassifier(n_neighbors=5))  
]  
  
# Loop over each target in your target_data dict  
for target, (X_train, X_test, y_train, y_test) in target_data.items():  
    print(f"\n==== Results for target: {target} ===")  
    results = []  
    for name, model in models:  
        model.fit(X_train, y_train)  
        acc = accuracy_score(y_test, model.predict(X_test))  
        results.append((name, acc))  
        print(f"{name}: {acc:.4f}") # quick print as it trains  
  
    # Create and display a sorted DataFrame for this target  
    df_results = pd.DataFrame(results, columns=["Model", "Accuracy"]).sort_values("Accuracy", ascending=False)  
    print("\nAccuracy Table:")  
    print(df_results.to_string(index=False))
```

```
--- Training Classical Models for Multiple Targets ---
```

```
==== Results for target: Running ===  
Logistic Regression: 0.8066  
Random Forest: 0.8215  
SVM (RBF Kernel): 0.8264  
Gradient Boosting: 0.8331  
K-Nearest Neighbors: 0.8165
```

```
Accuracy Table:
```

Model	Accuracy
Gradient Boosting	0.833058
SVM (RBF Kernel)	0.826446

```
Random Forest 0.821488
K-Nearest Neighbors 0.816529
Logistic Regression 0.806612

==== Results for target: Climbing ====
Logistic Regression: 1.0000
Random Forest: 1.0000
SVM (RBF Kernel): 0.9983
Gradient Boosting: 1.0000
K-Nearest Neighbors: 0.9884
```

Accuracy Table:

Model	Accuracy
Logistic Regression	1.000000
Random Forest	1.000000
Gradient Boosting	1.000000
SVM (RBF Kernel)	0.998347
K-Nearest Neighbors	0.988430

```
==== Results for target: Chasing ====
Logistic Regression: 1.0000
Random Forest: 1.0000
SVM (RBF Kernel): 1.0000
Gradient Boosting: 1.0000
K-Nearest Neighbors: 1.0000
```

Accuracy Table:

Model	Accuracy
Logistic Regression	1.0
Random Forest	1.0
SVM (RBF Kernel)	1.0
Gradient Boosting	1.0
K-Nearest Neighbors	1.0

▼ HA! More leakage!

It's a good thing to notice when accuracy is perfect.. we know that's garbage! (Or really bad data.)

So let's test this and see..

```
print("\n--- Squirrel Target Audit: Leakage, Balance, Duplication ---")

for target, (X_train, X_test, y_train, y_test) in target_data.items():
    print(f"\n🔍 Target: {target}")

    # Check if X_train is a DataFrame
    if isinstance(X_train, pd.DataFrame):
        if target in X_train.columns:
            print(f"⚠️ Target '{target}' found in feature set - potential leak")
```

```

        feature_count = X_train.shape[1]
    else:
        print("ℹ Feature matrix is a NumPy array – skipping column name check")
        feature_count = X_train.shape[1]

    # Class balance
    train_dist = pd.Series(y_train).value_counts(normalize=True).round(3)
    test_dist = pd.Series(y_test).value_counts(normalize=True).round(3)
    print("📊 Class distribution (Train):")
    print(train_dist.to_string())
    print("📊 Class distribution (Test):")
    print(test_dist.to_string())

    # Check for duplicate rows across train/test (only if DataFrame)
    if isinstance(X_train, pd.DataFrame) and isinstance(X_test, pd.DataFrame):
        overlap = pd.merge(X_train, X_test, how='inner')
        print(f"⌚ Duplicate rows between train/test: {len(overlap)}")
    else:
        print("⌚ Skipping duplicate check – feature matrices are arrays.")

    # Feature count
    print(f"🧠 Feature count: {feature_count}")

```

--- Squirrel Target Audit: Leakage, Balance, Duplication ---

```

🔍 Target: Running
ℹ Feature matrix is a NumPy array – skipping column name check.
📊 Class distribution (Train):
Running
0    0.758
1    0.242
📊 Class distribution (Test):
Running
0    0.759
1    0.241
⌚ Skipping duplicate check – feature matrices are arrays.
🧠 Feature count: 36

🔍 Target: Climbing
ℹ Feature matrix is a NumPy array – skipping column name check.
📊 Class distribution (Train):
Climbing
0    0.782
1    0.218
📊 Class distribution (Test):
Climbing
0    0.782
1    0.218
⌚ Skipping duplicate check – feature matrices are arrays.
🧠 Feature count: 36

```

```
🔍 Target: Chasing
ℹ️ Feature matrix is a NumPy array – skipping column name check.
📊 Class distribution (Train):
Chasing
0    0.908
1    0.092
📊 Class distribution (Test):
Chasing
0    0.907
1    0.093
🔄 Skipping duplicate check – feature matrices are arrays.
🧠 Feature count: 36
```

```
from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay, f1_score

print("\n--- Classical Model Evaluation with Confusion Matrix & F1 Score ---")

for target, (X_train, X_test, y_train, y_test) in target_data.items():
    print(f"\n🔍 Target: {target}")
    for name, model in models:
        model.fit(X_train, y_train)
        y_pred = model.predict(X_test)
        acc = accuracy_score(y_test, y_pred)
        f1 = f1_score(y_test, y_pred, average='binary')
        print(f"{name}: Accuracy = {acc:.4f}, F1 Score = {f1:.4f}")

        # Confusion matrix
        print("Confusion Matrix:")
        cm = confusion_matrix(y_test, y_pred)
        disp = ConfusionMatrixDisplay(confusion_matrix=cm)
        disp.plot(cmap='Blues')
        plt.title(f"{name} - {target}")
        plt.show()
```


▼ Stick with Running Squirrels

Both climbing and chasing are either too small to test/train with, or the data is too easy to separate.. and there's no need to test it.

For example, chasing is listed in around 9% percent of the data rows. So that's basically like a 9% percent chance of rain in the forecast.. and you can virtually guess the 91% portion with ease.

At least we tried to see what's possible. Or not.

```
# =====
# === FEATURE ENGINEERING: Domain-Driven Features =====
# =====

# 🧠 Activity Score: behavioral intensity
squirrel_df['activity_score'] = squirrel_df[['Running', 'Chasing', 'Eating', 'Fo

# 🧠 Vocalization Score: expressiveness
squirrel_df['vocalization_score'] = squirrel_df[['Kuks', 'Quaas', 'Moans']].asty|
```

```

# 🐿 Tail Signal Score: communication signals
# Use only columns that exist in the DataFrame
tail_cols = [col for col in ['Tail_flags', 'Tail_wags'] if col in squirrel_df.columns]
squirrel_df['tail_signal_score'] = squirrel_df[tail_cols].astype(int).sum(axis=1)

print("\n✅ Domain-driven features created:")
print(squirrel_df[['activity_score', 'vocalization_score', 'tail_signal_score']])

```

✅ Domain-driven features created:

	activity_score	vocalization_score	tail_signal_score
0	0	0	0.0
1	0	0	0.0
2	1	0	0.0
3	2	0	0.0
4	1	0	0.0

Domain-Driven Features

We created three behavioral features:

- **Activity Score:** Combines `running`, `chasing`, `eating`, and `foraging` to measure overall behavioral intensity.
- **Vocalization Score:** Aggregates `kuks`, `quaas`, and `moans` to reflect expressiveness. Vocal squirrels may be more reactive or social, and may have been more noticeable for this "census".
- **Tail Signal Score:** Combines `tail_flags` and `tail_wags` (when available) to capture non-verbal communication.

We can encode behavioral or communication elements that basic raw data or fur colors can't totally synthesize.

We'll test this expressive score after we try a basic run with 'Running'.

Feature Engineering and Final Preparation for the Squirrel Data

Alright team, time to get our squirrel dataset completely ready for building models! We've cleaned it up, and now we need to make sure all the information is in a format our models can understand and learn from.

Then we'll test it with a few classic models, like we did with the soccer dataset.

```
import pandas as pd
```

```

import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
from sklearn.svm import SVC
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score
from sklearn.preprocessing import OneHotEncoder, StandardScaler, LabelEncoder
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer
from xgboost import XGBClassifier

print("--- Starting Full Preprocessing and Classical Model Comparison for Squirrels ---")

# --- Step 1: Feature Engineering - Convert Boolean Columns to Integers ---
# This is a crucial step to make sure the boolean columns are numeric.
boolean_cols = squirrel_df.select_dtypes(include='bool').columns.tolist()

for col in boolean_cols:
    squirrel_df[col] = squirrel_df[col].astype(int)

print("Boolean columns successfully converted to integers.")

# --- Step 2: Define the target variable and features to use ---
# We use all relevant features, explicitly excluding identifier columns.
target_variable = 'Running'

# Define feature lists after boolean conversion
numerical_features = ["X", "Y", "Above Ground Sighter Measurement"]
categorical_features = ["Age", "Primary Fur Color", "Highlight Fur Color", "Location"]
behavioral_features = [
    'Chasing', 'Climbing', 'Eating', 'Foraging', 'Kuks', 'Quaas',
    'Moans', 'Tail flags', 'Tail twitches', 'Approaches', 'Indifferent',
    'Runs from'
]
]

# Separate features (X) and target variable (y)
X_squirrel = squirrel_df[numerical_features + categorical_features + behavioral_features]
y_squirrel = squirrel_df[target_variable].astype(int)

# --- Step 3: Preprocess Data with ColumnTransformer ---
# Use pipelines to handle both imputation and scaling for numerical data
numerical_pipeline = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='median')),
    ('scaler', StandardScaler())
])

preprocessor = ColumnTransformer(
    transformers=[
        ('num', numerical_pipeline, numerical_features + behavioral_features),
        ('cat', OneHotEncoder(), categorical_features)
    ]
)

```

```

        ],
        remainder='drop'
    )

# Apply preprocessing to all data
X_processed_squirrel = preprocessor.fit_transform(X_squirrel)

print("\nFeatures have been successfully encoded and scaled, ready for modeling.")
print(f"Final feature matrix shape: {X_processed_squirrel.shape}")

# --- Step 4: Split the Processed Data ---
X_train, X_test, y_train, y_test = train_test_split(
    X_processed_squirrel, y_squirrel, test_size=0.2, random_state=42, stratify=y
)

print("\nData successfully split into training and testing sets for modeling.")
print(f"Training features shape: {X_train.shape}")
print(f"Testing features shape: {X_test.shape}")

# --- Step 5: Define and Train Models ---
print("\n--- Training Classical Models for Squirrel Data ---")

# Encode labels for XGBoost, which requires integer targets
le = LabelEncoder()
y_train_enc = le.fit_transform(y_train)
y_test_enc = le.transform(y_test)

models = [
    ("Logistic Regression", LogisticRegression(max_iter=1000, random_state=42)),
    ("Random Forest", RandomForestClassifier(random_state=42)),
    ("SVM (RBF Kernel)", SVC(kernel='rbf', random_state=42)),
    ("Gradient Boosting", GradientBoostingClassifier(random_state=42)),
    ("K-Nearest Neighbors", KNeighborsClassifier(n_neighbors=5)),
    ("XGBoost", XGBClassifier(use_label_encoder=False, eval_metric='logloss', rai
]

results = []

for name, model in models:
    if name == "XGBoost":
        model.fit(X_train, y_train_enc)
        acc = accuracy_score(y_test_enc, model.predict(X_test))
    else:
        model.fit(X_train, y_train)
        acc = accuracy_score(y_test, model.predict(X_test))
    results.append((name, acc))

# --- Step 6: Create and Display Results Table ---
df_results = pd.DataFrame(results, columns=["Model", "Accuracy"]).sort_values(by:

```

```

print("\n== Classical Model Comparison for Squirrel Data ==")
print(df_results.to_string(index=False))

print("\n--- Model Testing Complete! 🎉 ---")

--- Starting Full Preprocessing and Classical Model Comparison for Squirrel Data
Boolean columns successfully converted to integers.

Features have been successfully encoded and scaled, ready for modeling.
Final feature matrix shape: (3023, 36)

Data successfully split into training and testing sets for modeling.
Training features shape: (2418, 36)
Testing features shape: (605, 36)

--- Training Classical Models for Squirrel Data ---

==== Classical Model Comparison for Squirrel Data ===
      Model Accuracy
Gradient Boosting  0.833058
      SVM (RBF Kernel)  0.826446
      Random Forest  0.821488
K-Nearest Neighbors  0.816529
      XGBoost  0.813223
Logistic Regression  0.806612

--- Model Testing Complete! 🎉 ---
/usr/local/lib/python3.12/dist-packages/xgboost/training.py:183: UserWarning: [1
Parameters: { "use_label_encoder" } are not used.

bst.update(dtrain, iteration=i, fobj=obj)

```

```

# =====
# === SIDE TEST: expressive_combo Feature Impact =====
# =====

import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
from sklearn.svm import SVC
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score
from sklearn.preprocessing import OneHotEncoder, StandardScaler, LabelEncoder
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer
from xgboost import XGBClassifier

print(" --- Starting Side Test with expressive_combo Feature ---")

```

```

# --- Step 1: Feature Engineering - Convert Boolean Columns to Integers ---
boolean_cols = squirrel_df.select_dtypes(include='bool').columns.tolist()
for col in boolean_cols:
    squirrel_df[col] = squirrel_df[col].astype(int)

# 🖌 Interaction Feature: vocal x tail signal
squirrel_df['expressive_combo'] = (
    squirrel_df[['Kuks', 'Quaas', 'Moans']].astype(int).sum(axis=1)
    * squirrel_df[['Tail flags', 'Tail twitches']].astype(int).sum(axis=1)
)

print("Boolean columns converted and expressive_combo feature created.")

# --- Step 2: Define Features and Target ---
target_variable = 'Running'

numerical_features_test = ["X", "Y", "Above Ground Sighter Measurement", "expres
categorical_features = ["Age", "Primary Fur Color", "Highlight Fur Color", "Loc
behavioral_features = [
    'Chasing', 'Climbing', 'Eating', 'Foraging', 'Kuks', 'Quaas',
    'Moans', 'Tail flags', 'Tail twitches', 'Approaches', 'Indifferent',
    'Runs from'
]

X_squirrel_test = squirrel_df[numerical_features_test + categorical_features +
y_squirrel_test = squirrel_df[target_variable].astype(int)

# --- Step 3: Preprocessing ---
numerical_pipeline = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='median')),
    ('scaler', StandardScaler())
])

preprocessor = ColumnTransformer(
    transformers=[
        ('num', numerical_pipeline, numerical_features_test + behavioral_featur
        ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features)
    ],
    remainder='drop'
)

X_processed_test = preprocessor.fit_transform(X_squirrel_test)

print("\nFeatures encoded and scaled.")
print(f"Processed feature matrix shape: {X_processed_test.shape}")

# --- Step 4: Train/Test Split ---
X_train_t, X_test_t, y_train_t, y_test_t = train_test_split(
    X_processed_test, y_squirrel_test, test_size=0.2, random_state=42, stratify
)

```

```

print("Data split complete.")
print(f"Train shape: {X_train_t.shape}, Test shape: {X_test_t.shape}")

# --- Step 5: Train Models ---
models = [
    ("Logistic Regression", LogisticRegression(max_iter=1000, random_state=42)),
    ("Random Forest", RandomForestClassifier(random_state=42)),
    ("SVM (RBF Kernel)", SVC(kernel='rbf', random_state=42)),
    ("Gradient Boosting", GradientBoostingClassifier(random_state=42)),
    ("K-Nearest Neighbors", KNeighborsClassifier(n_neighbors=5)),
    ("XGBoost", XGBClassifier(use_label_encoder=False, eval_metric='logloss', r
]

le = LabelEncoder()
y_train_enc = le.fit_transform(y_train_t)
y_test_enc = le.transform(y_test_t)

results_test = []

for name, model in models:
    if name == "XGBoost":
        model.fit(X_train_t, y_train_enc)
        acc = accuracy_score(y_test_enc, model.predict(X_test_t))
    else:
        model.fit(X_train_t, y_train_t)
        acc = accuracy_score(y_test_t, model.predict(X_test_t))
    results_test.append((name, acc))

# --- Step 6: Display Results ---
df_results_test = pd.DataFrame(results_test, columns=["Model", "Accuracy"]).sort_index()

print("\n==== Side Test Results with expressive_combo Feature ===")
print(df_results_test.to_string(index=False))

print("\n--- Side Test Complete! 🎯 ---")

--- Starting Side Test with expressive_combo Feature ---
Boolean columns converted and expressive_combo feature created.

Features encoded and scaled.
Processed feature matrix shape: (3023, 37)
Data split complete.
Train shape: (2418, 37), Test shape: (605, 37)
/usr/local/lib/python3.12/dist-packages/xgboost/training.py:183: UserWarning: [1
Parameters: { "use_label_encoder" } are not used.

bst.update(dtrain, iteration=i, fobj=obj)

==== Side Test Results with expressive_combo Feature ====
      Model  Accuracy

```

	Model Accuracy
Gradient Boosting	0.836364
Random Forest	0.829752
SVM (RBF Kernel)	0.824793
K-Nearest Neighbors	0.816529
XGBoost	0.813223
Logistic Regression	0.806612

--- Side Test Complete!  ---

```

import matplotlib.pyplot as plt
import pandas as pd

# Original and side test results
original_acc = {
    "Gradient Boosting": 0.833058,
    "SVM (RBF Kernel)": 0.826446,
    "Random Forest": 0.821488,
    "K-Nearest Neighbors": 0.816529,
    "XGBoost": 0.813223,
    "Logistic Regression": 0.806612
}

combo_acc = {
    "Gradient Boosting": 0.836364,
    "SVM (RBF Kernel)": 0.824793,
    "Random Forest": 0.829752,
    "K-Nearest Neighbors": 0.816529,
    "XGBoost": 0.813223,
    "Logistic Regression": 0.806612
}

# Create DataFrame
df_compare = pd.DataFrame({
    "Original Accuracy": original_acc,
    "With expressive_combo": combo_acc
}).sort_values(by="With expressive_combo", ascending=False)

# Plot
plt.figure(figsize=(10, 6))
bars1 = plt.bar(df_compare.index, df_compare["Original Accuracy"], color="#6baed6")
bars2 = plt.bar(df_compare.index, df_compare["With expressive_combo"], color="#f9c232")

# Annotate bars with delta
for i, model in enumerate(df_compare.index):
    delta = df_compare["With expressive_combo"][model] - df_compare["Original Accuracy"][model]
    plt.text(i, df_compare["With expressive_combo"][model] + 0.001, f"+{delta:.3f}",
             ha='center', va='bottom', fontsize=10, color='black')

plt.title("🎯 Model Accuracy Comparison: Original vs expressive_combo", fontsize=14)
plt.ylabel("Accuracy")

```

```
plt.ylim(0.80, 0.84)
plt.xticks(rotation=45)
plt.legend()
plt.grid(axis='y', linestyle='--', alpha=0.5)
plt.tight_layout()
plt.show()
```

So.. does that mean our original model is garbage?
NOPE!

It's more to do with the datasets.

Squirrels.. well, they run, they eat, they climb. And that's kinda it. There's not too much random or unpredictability with squirrels.

However, if we decided to use GB or RF with the expressive combo we'd have even more accuracy. *BUT* we're also combining two or three features into one, so it's somewhat clumping features together.. which helps boost the probability and prediction analysis.

```
#  
# SQUIRRLRRRRRRRRRRRRRRRRRRRRRRRRRRRRRRRRREL!  
#  
# --- Neural Network for Squirrel Data ---  
print("\n--- Starting Neural Network Training for Squirrel Data ---")  
  
from tensorflow.keras.models import Sequential  
from tensorflow.keras.layers import Dense  
from tensorflow.keras.optimizers import Adam  
from tensorflow.keras.utils import to_categorical  
from sklearn.model_selection import train_test_split  
import matplotlib.pyplot as plt  
  
# Use cleaned and preprocessed data  
X_train, X_test, y_train, y_test = train_test_split(  
    X_processed_squirrel, y_squirrel,  
    test_size=0.2, random_state=42, stratify=y_squirrel  
)  
  
# One-hot encode the target variables for softmax output  
y_train_encoded_squirrel = to_categorical(y_train)  
y_test_encoded_squirrel = to_categorical(y_test)  
  
# Build a shallow network (project requirement)  
input_dim = X_train.shape[1]  
output_dim = y_train_encoded_squirrel.shape[1]  
  
model_nn = Sequential([  
    Dense(32, activation='relu', input_shape=(input_dim,)),  
    Dense(16, activation='relu'),  
    Dense(output_dim, activation='softmax')  
)  
  
model_nn.compile(optimizer=Adam(learning_rate=0.001),  
                  loss='categorical_crossentropy',  
                  metrics=['accuracy'])  
  
model_nn.summary()  
  
# Train the model
```

```
... train the model

history_squirrel = model_nn.fit(
    X_train,
    y_train_encoded_squirrel,
    epochs=20,
    batch_size=32,
    validation_data=(X_test, y_test_encoded_squirrel),
    verbose=1
)

# Evaluate on test data
loss_nn, acc_nn = model_nn.evaluate(X_test, y_test_encoded_squirrel, verbose=0)
print(f"\n\ufe0f Neural Network Test Accuracy (Squirrel Data): {acc_nn:.4f}")

# Plot training history
plt.figure(figsize=(6,4))
plt.plot(history_squirrel.history['accuracy'], label='Train Acc')
plt.plot(history_squirrel.history['val_accuracy'], label='Val Acc')
plt.xlabel('Epoch')
plt.ylabel('Accuracy')
plt.title('Squirrel NN Accuracy')
plt.legend()
plt.show()

plt.figure(figsize=(6,4))
plt.plot(history_squirrel.history['loss'], label='Train Loss')
plt.plot(history_squirrel.history['val_loss'], label='Val Loss')
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.title('Squirrel NN Loss')
plt.legend()
plt.show()
```


- ▼ Not bad!

As we can see, the accuracy of the NN increased once it started to process the data again and again.

Keeping the outliers certainly seemed to help round out and assist the model, as well as dealing with the *NaN* identifiers.

In fact.. **why don't we test that out before we continue?!**

Testing GOOD v BAD Pipelines:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
```

```

from sklearn.model_selection import train_test_split
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense
from tensorflow.keras.optimizers import Adam

# --- GOOD PIPELINE ---
print("\n--- GOOD PIPELINE ---")

X_train_good, X_test_good, y_train_good, y_test_good = train_test_split(
    X_processed_squirrel, y_squirrel,
    test_size=0.2, random_state=42, stratify=y_squirrel
)

model_good = Sequential([
    tf.keras.Input(shape=(X_train_good.shape[1],)),
    Dense(32, activation='relu'),
    Dense(16, activation='relu'),
    Dense(2, activation='softmax')
])

model_good.compile(optimizer=Adam(learning_rate=0.001),
                    loss='sparse_categorical_crossentropy',
                    metrics=['accuracy'])

model_good.fit(X_train_good, y_train_good,
               validation_data=(X_test_good, y_test_good),
               epochs=20, batch_size=32, verbose=0)

loss_good, acc_good = model_good.evaluate(X_test_good, y_test_good, verbose=0)
print(f"✅ GOOD Pipeline Test Accuracy: {acc_good:.4f}")

# --- BAD PIPELINE ---
print("\n--- BAD PIPELINE ---")

bad_df = squirrel_df.copy()

# Remove outliers from AGSM
col = "Above Ground Sighter Measurement"
Q1 = bad_df[col].quantile(0.25)
Q3 = bad_df[col].quantile(0.75)
IQR = Q3 - Q1
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR
bad_df = bad_df[(bad_df[col] >= lower_bound) & (bad_df[col] <= upper_bound)]
print(f"Outliers removed. Rows left: {bad_df.shape[0]}")

# Define features and target

```

```

X_bad = bad_df.drop(columns=['Running'])
y_bad = bad_df['Running'].astype(int)

# --- Correlation check for leakage ---
print("\n🔍 Checking numeric feature correlations with 'Running':")
corrs = bad_df.corr(numeric_only=True)[['Running']].sort_values(ascending=False)
print(corrs)
perfect_corrs = corrs[abs(corrs) == 1].drop('Running', errors='ignore')
if not perfect_corrs.empty:
    print("\n⚠️ Perfect correlation(s) found - possible leakage:")
    print(perfect_corrs)

# Identify raw feature types
numerical_features_bad = X_bad.select_dtypes(include='number').columns.tolist()
categorical_features_bad = X_bad.select_dtypes(include=['object', 'category']).columns.tolist()

# BAD preprocessing (no imputation)
numerical_pipeline_bad = Pipeline(steps=[
    ('scaler', StandardScaler())
])
categorical_pipeline_bad = Pipeline(steps=[
    ('encoder', OneHotEncoder(handle_unknown='ignore'))
])

preprocessor_bad = ColumnTransformer(
    transformers=[
        ('num', numerical_pipeline_bad, numerical_features_bad),
        ('cat', categorical_pipeline_bad, categorical_features_bad)
    ],
    remainder='drop'
)

# --- SPLIT FIRST to avoid leakage ---
X_train_bad, X_test_bad, y_train_bad, y_test_bad = train_test_split(
    X_bad, y_bad, test_size=0.2, random_state=42, stratify=y_bad
)

# Fit preprocessor only on training data
X_train_bad = preprocessor_bad.fit_transform(X_train_bad)
X_test_bad = preprocessor_bad.transform(X_test_bad)

print("Preprocessing successful. Shapes:",
      X_train_bad.shape, X_test_bad.shape)

# Build and train NN
model_bad = Sequential([
    tf.keras.Input(shape=(X_train_bad.shape[1],)),
    Dense(32, activation='relu'),
    Dense(16, activation='relu'),
    Dense(2, activation='softmax')
])

```

```

// 

model_bad.compile(optimizer=Adam(learning_rate=0.001),
                   loss='sparse_categorical_crossentropy',
                   metrics=['accuracy'])

model_bad.fit(X_train_bad, y_train_bad,
               validation_data=(X_test_bad, y_test_bad),
               epochs=20, batch_size=32, verbose=0)

loss_bad, acc_bad = model_bad.evaluate(X_test_bad, y_test_bad, verbose=0)
print(f"❌ BAD Pipeline Test Accuracy: {acc_bad:.4f}")

# --- Summary Table ---
summary_df = pd.DataFrame({
    'Pipeline': ['GOOD (Imputed + Engineered)', 'BAD (NaNs + No Outliers)'],
    'Description': ['Cleaned + engineered features', 'Raw features, no fixes'],
    'Test Accuracy': [acc_good, acc_bad]
})

print("\n== COMPARISON TABLE ==")
print(summary_df.to_string(index=False))

# --- Bar Chart ---
plt.style.use('seaborn-v0_8')
colors = ['#4CAF50', '#F44336']

fig, ax = plt.subplots(figsize=(8, 6))
bars = ax.bar(summary_df['Pipeline'], summary_df['Test Accuracy'], color=colors

for bar in bars:
    height = bar.get_height()
    ax.text(bar.get_x() + bar.get_width()/2, height + 0.005, f'{height:.4f}',
            ha='center', va='bottom', fontsize=12)

ax.set_title('Neural Network Accuracy: GOOD vs BAD Pipeline', fontsize=14)
ax.set_ylabel('Test Accuracy', fontsize=12)
ax.set_ylim(min(acc_bad, acc_good) - 0.02, max(acc_bad, acc_good) + 0.02)
plt.tight_layout()
plt.show()

```


Oooooh.. even leakage in a "bad" pipeline test?

Alright.. let's fix this as well - maybe we'll just drop 'Running' (or if there's a duplicate floating somewhere).

After testing for both, there's something else going on.

```
# Count how many rows have a unique feature combination
feature_df = X_bad.copy()
unique_patterns = feature_df.drop_duplicates().shape[0]
total_rows = feature_df.shape[0]
print(f"Unique feature patterns: {unique_patterns} / {total_rows}")
```

```
Unique feature patterns: 2346 / 2346
```

Ah - bingo! Of course it's 1.000 accuracy if everything is unique!

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense
from tensorflow.keras.optimizers import Adam

print("\n--- BAD PIPELINE (Aggressive De-ID & Coarsening) ---")

# Copy raw data
bad_df = squirrel_df.copy()

# Remove outliers from AGSM
col = "Above Ground Sighter Measurement"
```

```

Q1 = bad_df[col].quantile(0.25)
Q3 = bad_df[col].quantile(0.75)
IQR = Q3 - Q1
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR
bad_df = bad_df[(bad_df[col] >= lower_bound) & (bad_df[col] <= upper_bound)]
print(f"Outliers removed. Rows left: {bad_df.shape[0]}")

# Drop obvious identifiers
id_like_cols = ['Unique Squirrel ID', 'Lat/Long', 'Date', 'Hectare Squirrel Number']
bad_df = bad_df.drop(columns=id_like_cols, errors='ignore')

# Aggressive binning for numeric coords
if 'X' in bad_df.columns:
    bad_df['X'] = (bad_df['X'] // 100) * 100
if 'Y' in bad_df.columns:
    bad_df['Y'] = (bad_df['Y'] // 100) * 100

# Group rare categories
categorical_cols = bad_df.select_dtypes(include=['object', 'category']).columns
for col in categorical_cols:
    freqs = bad_df[col].value_counts(normalize=True)
    rare_cats = freqs[freqs < 0.05].index
    bad_df[col] = bad_df[col].replace(rare_cats, 'Other')

# Define features and target
y_bad = bad_df['Running'].astype(int)
X_bad = bad_df.drop(columns=['Running'])

# Check unique patterns after aggressive de-ID
unique_patterns = X_bad.drop_duplicates().shape[0]
print(f"Unique feature patterns after aggressive de-ID: {unique_patterns} / {X_bad.shape[0]}")

# Identify raw feature types
numerical_features_bad = X_bad.select_dtypes(include='number').columns.tolist()
categorical_features_bad = X_bad.select_dtypes(include=['object', 'category']).columns.tolist()

# BAD preprocessing (no imputation)
numerical_pipeline_bad = Pipeline(steps=[
    ('scaler', StandardScaler())
])
categorical_pipeline_bad = Pipeline(steps=[
    ('encoder', OneHotEncoder(handle_unknown='ignore'))
])

preprocessor_bad = ColumnTransformer(
    transformers=[
        ('num', numerical_pipeline_bad, numerical_features_bad),
        ('cat', categorical_pipeline_bad, categorical_features_bad)
    ],
    remainder='drop'
)

```

```

        remainder= arrop
    )

# Split first to avoid leakage
X_train_bad, X_test_bad, y_train_bad, y_test_bad = train_test_split(
    X_bad, y_bad, test_size=0.2, random_state=42, stratify=y_bad
)

# Fit preprocessor only on training data
X_train_bad = preprocessor_bad.fit_transform(X_train_bad)
X_test_bad = preprocessor_bad.transform(X_test_bad)

# Build and train NN
model_bad = Sequential([
    tf.keras.Input(shape=(X_train_bad.shape[1],)),
    Dense(32, activation='relu'),
    Dense(16, activation='relu'),
    Dense(2, activation='softmax')
])
model_bad.compile(optimizer=Adam(learning_rate=0.001),
                   loss='sparse_categorical_crossentropy',
                   metrics=['accuracy'])

model_bad.fit(X_train_bad, y_train_bad,
              validation_data=(X_test_bad, y_test_bad),
              epochs=20, batch_size=32, verbose=0)

loss_bad, acc_bad = model_bad.evaluate(X_test_bad, y_test_bad, verbose=0)
print(f"🟡 BAD Pipeline (Aggressive De-ID) Test Accuracy: {acc_bad:.4f}")

# --- Bar Chart comparing GOOD vs BAD (Aggressive De-ID) ---
summary_df = pd.DataFrame({
    'Pipeline': ['GOOD (Imputed + Engineered)', 'BAD (Aggressive De-ID)'],
    'Description': ['Cleaned + engineered features', 'Identifiers removed, coar'],
    'Test Accuracy': [acc_good, acc_bad]
})

plt.style.use('seaborn-v0_8')
colors = ['#4CAF50', '#F44336']

fig, ax = plt.subplots(figsize=(8, 6))
bars = ax.bar(summary_df['Pipeline'], summary_df['Test Accuracy'], color=colors

for bar in bars:
    height = bar.get_height()
    ax.text(bar.get_x() + bar.get_width()/2, height + 0.005, f'{height:.4f}',
            ha='center', va='bottom', fontsize=12)

ax.set_title('Neural Network Accuracy: GOOD vs BAD Pipeline (Aggressive De-ID)')
ax.set_ylabel('Test Accuracy', fontsize=12)

```

```
ax.set_ylim(min(acc_bad, acc_good) - 0.02, max(acc_bad, acc_good) + 0.02)
plt.tight_layout()
plt.show()
```

NOT WHAT IT WAS - There *was* a 7% lower number for
BAD..

But now the model is going loopy. Internally. And this will be just another detour. So let's just talk about the prior one..

It may not be huge, but it's significant. Since it shows that, *as a baseline*, the data shape works a lot better with outliers and no NaN issues. There's some wiggle room because we're not using specific numbers for the test/train split, so unless we use a constant "random_state" number, we'll get fractional differences.

For now, the model realizes if a squirrel is 80 feet in the air, it's probably climbing, not running. Maybe chasing, maybe eating.. but we can tweak this from a pretty solid base.

Bottom line: the performance gap is meaningful.

```
# =====
# ===== PART 3: BASELINE & KERNEL METHODS =====
# =====

import time
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.pipeline import Pipeline
from sklearn.metrics import mean_squared_error
from sklearn.linear_model import LinearRegression, Ridge, Lasso, ElasticNet
from sklearn.kernel_ridge import KernelRidge
from sklearn.compose import ColumnTransformer

# --- Project Setup ---
RANDOM_STATE = 42
np.random.seed(RANDOM_STATE)

# --- Load NYC Squirrel Dataset ---
data_url = 'https://data.cityofnewyork.us/api/views/vfnx-vebw/rows.csv?accessType=Read&method=Get'
df = pd.read_csv(data_url)

# --- Clean & Prepare ---
# Fix column names for consistency
df.columns = df.columns.str.strip().str.replace(' ', '_').str.lower()

# Impute missing values
df['age'] = df['age'].fillna('unknown')
df['primary_fur_color'] = df['primary_fur_color'].fillna('gray')
df['approaches'] = df['approaches'].fillna(False).astype(int)

# Engineer features
df['activity_score'] = df[['running', 'chasing', 'eating', 'foraging']].astype(bool).sum(1)
```

```

# Define features and target
numerical_features = ['x', 'y', 'activity_score']
categorical_features = ['primary_fur_color', 'age']
target_col = 'approaches'

X = df[numerical_features + categorical_features]
y = df[target_col]

# --- Preprocessing Pipeline ---
preprocessor = ColumnTransformer(transformers=[
    ('num', StandardScaler(), numerical_features),
    ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features)
])

# --- Train/Test Split ---
X_train, X_val, y_train, y_val = train_test_split(
    X, y, test_size=0.2, random_state=RANDOM_STATE
)

# --- Helper Function ---
def fit_time_mse(pipeline):
    t0 = time.perf_counter()
    pipeline.fit(X_train, y_train)
    t1 = time.perf_counter()
    y_pred = pipeline.predict(X_val)
    mse = mean_squared_error(y_val, y_pred)
    return round(t1 - t0, 4), round(mse, 4)

print("\n--- Baseline Regressors ---")
results = []

# OLS
ols = Pipeline([('prep', preprocessor), ('model', LinearRegression())])
t, mse = fit_time_mse(ols)
results.append({'Model': 'OLS', 'Params': '—', 'Time': t, 'MSE': mse})
print(f"OLS → Time: {t}s | MSE: {mse}")

# Ridge
for alpha in [0.01, 0.1, 1, 10]:
    ridge = Pipeline([('prep', preprocessor), ('model', Ridge(alpha=alpha))])
    t, mse = fit_time_mse(ridge)
    results.append({'Model': 'Ridge', 'Params': f'alpha={alpha}', 'Time': t, 'MSE': mse})

# LASSO
for alpha in [0.01, 0.1, 1]:
    lasso = Pipeline([('prep', preprocessor), ('model', Lasso(alpha=alpha))])
    t, mse = fit_time_mse(lasso)
    results.append({'Model': 'LASSO', 'Params': f'alpha={alpha}', 'Time': t, 'MSE': mse})

```

```

# ElasticNet
for alpha in [0.01, 0.1]:
    for l1 in [0.2, 0.5, 0.8]:
        enet = Pipeline([('prep', preprocessor), ('model', ElasticNet(alpha=alpha, l1=l1))])
        t, mse = fit_time_mse(enet)
        results.append({'Model': 'ElasticNet', 'Params': f'alpha={alpha}, l1={l1}'})

```

Baseline Regressors ---

OLS → Time: 0.0869s | MSE: 0.0487

Our **Ordinary Least Squares** result — ~0.03 (or 0.08, some runs) seconds to train and a mean squared error (MSE) of 0.0487 on the validation set.

So this is our baseline.

```

print("\n--- Kernel Ridge Grid Search ---")
alphas = [0.01, 0.1, 1]
gammas = [0.01, 0.1, 1]
degrees = [2, 3]

# RBF Kernel
for alpha in alphas:
    for gamma in gammas:
        kr_rbf = Pipeline([
            ('prep', preprocessor),
            ('model', KernelRidge(kernel='rbf', alpha=alpha, gamma=gamma))
        ])
        t, mse = fit_time_mse(kr_rbf)
        results.append({'Model': 'KernelRidge-RBF', 'Params': f'alpha={alpha}, gamma={gamma}'})

# Polynomial Kernel
for alpha in alphas:
    for degree in degrees:
        kr_poly = Pipeline([
            ('prep', preprocessor),
            ('model', KernelRidge(kernel='polynomial', alpha=alpha, degree=degree))
        ])
        t, mse = fit_time_mse(kr_poly)
        results.append({'Model': 'KernelRidge-Poly', 'Params': f'alpha={alpha}, degree={degree}'})

```

Kernel Ridge Grid Search ---

And now for the results..

```
results_df = pd.DataFrame(results).sort_values(by='MSE')
print("\n--- Top Models by MSE ---")
print(results_df.head(10).to_string(index=False))

# Optional: plot MSEs
plt.figure(figsize=(10, 6))
sns.barplot(data=results_df, x='Model', y='MSE', hue='Params', dodge=False)
plt.xticks(rotation=45)
plt.title('Model Comparison: MSE on Approaches Prediction')
plt.tight_layout()
plt.show()
```

Hmmm.. helpful. Obviously the KR - RBF model is the worst one.

But let's see what else..

```
# =====
# === LINE PLOT: Ridge Regression MSE vs Alpha =====
# =====
import seaborn as sns
import matplotlib.pyplot as plt

print("\n--- Generating Line Plot for Ridge Regression ---")

# Filter Ridge results
ridge_df = results_df[results_df['Model'] == 'Ridge'].copy()

# Extract alpha from params string
ridge_df['alpha'] = ridge_df['Params'].str.extract(r'alpha=(\d+\.\?\d*)').astype(float)

print("Ridge Regression Results:")
print(ridge_df[['alpha', 'MSE']])

# Plot
plt.figure(figsize=(8, 5))
sns.lineplot(data=ridge_df, x='alpha', y='MSE', marker='o', linewidth=2)
plt.title('Ridge Regression: MSE vs Alpha')
plt.xlabel('Alpha (Regularization Strength)')
plt.ylabel('Mean Squared Error')
plt.grid(True)
plt.tight_layout()
plt.show()
```

So.. alpha is useless, basically. Looks like we fitted and scaled the data pretty well.

```
# =====
# === HEATMAP: Kernel Ridge RBF MSE by Alpha & Gamma =====
# =====
print("\n--- Generating Heatmap for Kernel Ridge (RBF) ---")

# Filter RBF results
rbf_df = results_df[results_df['Model'] == 'KernelRidge-RBF'].copy()

# Extract alpha and gamma
rbf_df['alpha'] = rbf_df['Params'].str.extract(r'alpha=(\d+\.\?\d*)').astype(float)
rbf_df['gamma'] = rbf_df['Params'].str.extract(r'gamma=(\d+\.\?\d*)').astype(float)

print("RBF Kernel Grid Search Results:")
print(rbf_df[['alpha', 'gamma', 'MSE']])

# Pivot for heatmap
heatmap_data = rbf_df.pivot(index='alpha', columns='gamma', values='MSE')
```

```
# Plot
plt.figure(figsize=(8, 6))
sns.heatmap(heatmap_data, annot=True, fmt=".4f", cmap='coolwarm')
plt.title('Kernel Ridge (RBF): MSE Heatmap')
plt.xlabel('Gamma')
plt.ylabel('Alpha')
plt.tight_layout()
plt.show()
print("\n=====")
print("🔥 HEATMAP: Kernel Ridge (RBF) – MSE by alpha & gamma 🔥")
print("=====")
print("\nColor Legend:")
print("- ⚪ Dark Blue = Low MSE (better)")
print("- ⚫ Light Blue = Slightly worse")
print("- ⚩ Yellowish = Moderate error")
print("- ⚫ Red = High MSE (worst)")
```

So a slightly increased alpha and virtually no gamma is optimal.

However, the really low alpha plus the high gamma equals (as Gemini AI stated well) a model that's "trying too hard to fit every bump in the data."

```
# =====
# === SCATTER PLOT: MSE vs Training Time =====
# =====
print("\n--- Generating Scatter Plot for MSE vs Training Time ---")

print("Model Performance Overview:")
print(results_df[['Model', 'Params', 'Time', 'MSE']].head(10))

plt.figure(figsize=(8, 6))
sns.scatterplot(data=results_df, x='Time', y='MSE', hue='Model', style='Model',
plt.title('Model Tradeoff: MSE vs Training Time')
plt.xlabel('Training Time (seconds)')
plt.ylabel('Validation MSE')
plt.grid(True)
plt.tight_layout()
plt.show()
```

That bottom left corner has a couple of really good options.. but our most optimal choice is listed below..

```
# =====
# === BEST MODEL SUMMARY =====
# =====
print("\n--- Best Model by MSE ---")
best_model = results_df.sort_values(by='MSE').iloc[0]
```

```
print(f"✅ Best Model: {best_model['Model']} with {best_model['Params']}")  
print(f"⌚ Training Time: {best_model['Time']} seconds")  
print(f"〽️ MSE: {best_model['MSE']}")
```

```
--- Best Model by MSE ---  
✅ Best Model: Ridge with alpha=10  
⌚ Training Time: 0.0123 seconds  
〽️ MSE: 0.0486
```

```
# ======  
# 🧠 SHALLOW NETWORK TRAINING – FULL PIPELINE (NaN-Safe, Numeric-Only)  
# ======  
  
import pandas as pd  
import numpy as np  
import torch  
import torch.nn as nn  
from torch.utils.data import TensorDataset, DataLoader  
from sklearn.preprocessing import StandardScaler  
import matplotlib.pyplot as plt  
  
# -----  
# 1. CLEAN & PREPARE DATA  
# -----  
print("NaNs in X_train before cleaning:", X_train.isna().sum().sum())  
print("NaNs in X_val before cleaning:", X_val.isna().sum().sum())  
  
# Drop columns that are entirely NaN  
X_train = X_train.dropna(axis=1, how='all')  
X_val = X_val.dropna(axis=1, how='all')  
  
# Keep only numeric columns  
X_train = X_train.select_dtypes(include=[np.number])  
X_val = X_val.select_dtypes(include=[np.number])  
  
# Fill remaining NaNs with training-set means  
train_means = X_train.mean()  
X_train = X_train.fillna(train_means)  
X_val = X_val.fillna(train_means) # use train means for val set  
  
# Fill NaNs in targets (regression)  
y_train = y_train.fillna(y_train.mean())  
y_val = y_val.fillna(y_val.mean())  
  
print("NaNs in X_train after cleaning:", X_train.isna().sum().sum())  
print("NaNs in X_val after cleaning:", X_val.isna().sum().sum())  
print("Any all-NaN columns left?", X_train.isna().all().any())  
#
```

```

# -----
# 2. SCALE FEATURES & TARGETS
# -----
input_size = X_train.shape[1]
output_size = 1 # Regression
print(f"✓ Input size: {input_size}, Output size: {output_size}")

scaler_X = StandardScaler()
X_train_scaled = scaler_X.fit_transform(X_train)
X_val_scaled = scaler_X.transform(X_val)

scaler_y = StandardScaler()
y_train_scaled = scaler_y.fit_transform(y_train.values.reshape(-1, 1))
y_val_scaled = scaler_y.transform(y_val.values.reshape(-1, 1))

print("Any NaNs in X_train_scaled?", np.isnan(X_train_scaled).any())
print("Any NaNs in y_train_scaled?", np.isnan(y_train_scaled).any())

# -----
# 3. CREATE TENSORS & DATALOADERS
# -----
X_train_tensor = torch.tensor(X_train_scaled, dtype=torch.float32)
y_train_tensor = torch.tensor(y_train_scaled, dtype=torch.float32)
X_val_tensor = torch.tensor(X_val_scaled, dtype=torch.float32)
y_val_tensor = torch.tensor(y_val_scaled, dtype=torch.float32)

train_dataset = TensorDataset(X_train_tensor, y_train_tensor)
val_dataset = TensorDataset(X_val_tensor, y_val_tensor)

train_loader = DataLoader(train_dataset, batch_size=64, shuffle=True)
val_loader = DataLoader(val_dataset, batch_size=64, shuffle=False)

print("✓ DataLoaders created successfully.")

# -----
# 4. ACTIVATION FUNCTION
# -----
class Swish(nn.Module):
    def forward(self, x):
        return x * torch.sigmoid(x)

# -----
# 5. SHALLOW NETWORK ARCHITECTURE
# -----
class ShallowNet(nn.Module):
    def __init__(self, input_size, hidden_size, output_size, enhanced=False):
        super().__init__()
        self.fc1 = nn.Linear(input_size, hidden_size)
        self.bn1 = nn.BatchNorm1d(hidden_size) if enhanced else nn.Identity()
        self.dropout = nn.Dropout(0.3) if enhanced else nn.Identity()
        self.activation = Swish() if enhanced else nn.ReLU()

```

```

        self.fc2 = nn.Linear(hidden_size, output_size)

    def forward(self, x):
        x = self.fc1(x)
        x = self.bn1(x)
        x = self.activation(x)
        x = self.dropout(x)
        return self.fc2(x)

# -----
# 6. TRAINING LOOP
# -----
def train_model(model, train_loader, val_loader, epochs=50, lr=1e-3):
    device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
    model.to(device)
    optimizer = torch.optim.Adam(model.parameters(), lr=lr)
    criterion = nn.MSELoss()
    train_losses, val_losses = [], []

    for epoch in range(epochs):
        model.train()
        train_loss = 0
        for xb, yb in train_loader:
            xb, yb = xb.to(device), yb.to(device)
            optimizer.zero_grad()
            loss = criterion(model(xb), yb)
            loss.backward()
            optimizer.step()
            train_loss += loss.item()
        train_losses.append(train_loss / len(train_loader))

        model.eval()
        val_loss = 0
        with torch.no_grad():
            for xb, yb in val_loader:
                xb, yb = xb.to(device), yb.to(device)
                val_loss += criterion(model(xb), yb).item()
        val_losses.append(val_loss / len(val_loader))

        if (epoch + 1) % 10 == 0 or epoch == 0:
            print(f"Epoch {epoch+1}/{epochs} - Train Loss: {train_losses[-1]:.4f}")

    return train_losses, val_losses, val_losses[-1]

# -----
# 7. PLOTTING HELPER
# -----
def plot_losses(train, val, title):
    plt.figure(figsize=(4, 3))
    plt.plot(train, label='Train')
    plt.plot(val, label='Val')

```

```
plt.plot(val, train - val, )
plt.title(title)
plt.xlabel("Epochs")
plt.ylabel("Loss")
plt.legend()
plt.tight_layout()
plt.show()

# -----
# 8. RUN BASELINE
# -----
print("\n◆ Training ShallowNet (Baseline)")
shallow_base = ShallowNet(input_size, 64, output_size, enhanced=False)
train_b, val_b, mse_b = train_model(shallow_base, train_loader, val_loader)
plot_losses(train_b, val_b, "ShallowNet Baseline")
print(f"▣ Final MSE (Baseline): {mse_b:.4f}")

# -----
# 9. RUN ENHANCED
# -----
print("\n◆ Training ShallowNet (Enhanced)")
shallow_enh = ShallowNet(input_size, 64, output_size, enhanced=True)
train_e, val_e, mse_e = train_model(shallow_enh, train_loader, val_loader)
plot_losses(train_e, val_e, "ShallowNet Enhanced")
print(f"▣ Final MSE (Enhanced): {mse_e:.4f}")
```

Even with just 3 features, it's clear that the baseline is a fraction slower or "worse" than the enhanced model.

Actually, unless other elements get introduced, it begs the question whether it's worth the extra processing power for the Enhanced -- when it comes to a dataset and architecture using this.

```
# =====
# 🧠 DEEP NETWORK (4 Hidden Layers) - NaN-Safe
# =====

class DeepNet(nn.Module):
    def __init__(self, input_size, hidden_size, output_size, enhanced=False):
        super().__init__()
        layers = []
        for i in range(11):
            if i < 4:
                layers.append(nn.Linear(input_size, hidden_size))
                layers.append(nn.ReLU())
                input_size = hidden_size
            else:
                layers.append(nn.Linear(hidden_size, output_size))
                layers.append(nn.ReLU())
                break
```

```
        for i in range(4):
            layers.append(nn.Linear(input_size if i == 0 else hidden_size, hidden_size))
            if enhanced:
                layers.append(nn.BatchNorm1d(hidden_size))
                layers.append(Swish())
                layers.append(nn.Dropout(0.3))
            else:
                layers.append(nn.ReLU())
        layers.append(nn.Linear(hidden_size, output_size))
        self.net = nn.Sequential(*layers)

    def forward(self, x):
        return self.net(x)

# --- Baseline ---
print("\n◆ Training DeepNet (Baseline)")
deep_base = DeepNet(input_size, 128, output_size, enhanced=False)
train_b, val_b, mse_b = train_model(deep_base, train_loader, val_loader)
plot_losses(train_b, val_b, "DeepNet Baseline")
print(f"▣ Final MSE (Baseline): {mse_b:.4f}")

# --- Enhanced ---
print("\n◆ Training DeepNet (Enhanced)")
deep_enh = DeepNet(input_size, 128, output_size, enhanced=True)
train_e, val_e, mse_e = train_model(deep_enh, train_loader, val_loader)
plot_losses(train_e, val_e, "DeepNet Enhanced")
print(f"▣ Final MSE (Enhanced): {mse_e:.4f}")
```

Hmmm.. looks like - other than the MSE difference - the Baseline started to do worse the longer it ran!

By coming closer, the model started to overfit.. not a good thing.

The Enhanced model started weird, but leveled out and stayed steady - thanks to some added elements that were like guardrails. Or like those rails in the gutter at the bowling alley.

```
# =====
# 🧠 WIDE NETWORK (Single Wide Layer) - NaN-Safe
# =====

class WideNet(nn.Module):
    def __init__(self, input_size, wide_size, output_size, enhanced=False):
        super().__init__()
        self.fc1 = nn.Linear(input_size, wide_size)
        self.bn1 = nn.BatchNorm1d(wide_size) if enhanced else nn.Identity()
        self.dropout = nn.Dropout(0.3) if enhanced else nn.Identity()
        self.activation = Swish() if enhanced else nn.ReLU()
        self.fc2 = nn.Linear(wide_size, output_size)
```

```
def forward(self, x):
    x = self.fc1(x)
    x = self.bn1(x)
    x = self.activation(x)
    x = self.dropout(x)
    return self.fc2(x)

# --- Baseline ---
print("\n◆ Training WideNet (Baseline)")
wide_base = WideNet(input_size, 512, output_size, enhanced=False)
train_b, val_b, mse_b = train_model(wide_base, train_loader, val_loader)
plot_losses(train_b, val_b, "WideNet Baseline")
print(f"▣ Final MSE (Baseline): {mse_b:.4f}")

# --- Enhanced ---
print("\n◆ Training WideNet (Enhanced)")
wide_enh = WideNet(input_size, 512, output_size, enhanced=True)
train_e, val_e, mse_e = train_model(wide_enh, train_loader, val_loader)
plot_losses(train_e, val_e, "WideNet Enhanced")
print(f"▣ Final MSE (Enhanced): {mse_e:.4f}")
```

OK.. so while this wasn't as bad of a difference between the two (as it was with the DeepNet structure), it's still pointing out that the basic ShallowNet model seems to be a good choice with decent MSE levels. For this dataset, or a similar one, at least. The DeepNet Enhanced model is just a minuscule amount less, but how much computing power would be needed compared to a simple ShallowNet model?

```
# =====
# [!] Ranked Model Comparison: MSE + Compute Cost Heatmaps
# =====

import torch
import torch.nn as nn
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# --- Swish activation ---
class Swish(nn.Module):
    def forward(self, x):
        return x * torch.sigmoid(x)

# --- Architectures ---
class ShallowNet(nn.Module):
    def __init__(self, input_size, hidden_size, output_size, enhanced=False):
        super().__init__()
        self.fc1 = nn.Linear(input_size, hidden_size)
        self.bn1 = nn.BatchNorm1d(hidden_size) if enhanced else nn.Identity()
        self.dropout = nn.Dropout(0.3) if enhanced else nn.Identity()
```

```

        self.activation = Swish() if enhanced else nn.ReLU()
        self.fc2 = nn.Linear(hidden_size, output_size)
    def forward(self, x):
        x = self.fc1(x); x = self.bn1(x); x = self.activation(x); x = self.dropout(x)
        return self.fc2(x)

class DeepNet(nn.Module):
    def __init__(self, input_size, hidden_size, output_size, enhanced=False):
        super().__init__()
        layers = []
        for i in range(4):
            layers.append(nn.Linear(input_size if i == 0 else hidden_size, hidden_size))
            if enhanced:
                layers.append(nn.BatchNorm1d(hidden_size))
                layers.append(Swish())
                layers.append(nn.Dropout(0.3))
            else:
                layers.append(nn.ReLU())
        layers.append(nn.Linear(hidden_size, output_size))
        self.net = nn.Sequential(*layers)
    def forward(self, x):
        return self.net(x)

class WideNet(nn.Module):
    def __init__(self, input_size, wide_size, output_size, enhanced=False):
        super().__init__()
        self.fc1 = nn.Linear(input_size, wide_size)
        self.bn1 = nn.BatchNorm1d(wide_size) if enhanced else nn.Identity()
        self.dropout = nn.Dropout(0.3) if enhanced else nn.Identity()
        self.activation = Swish() if enhanced else nn.ReLU()
        self.fc2 = nn.Linear(wide_size, output_size)
    def forward(self, x):
        x = self.fc1(x); x = self.bn1(x); x = self.activation(x); x = self.dropout(x)
        return self.fc2(x)

# --- Config ---
input_size = 3
output_size = 1

# --- Final MSEs from your runs ---
mse_values = {
    'ShallowNet Baseline': 0.8757,
    'ShallowNet Enhanced': 0.8744,
    'DeepNet Baseline': 0.9034,
    'DeepNet Enhanced': 0.8741,
    'WideNet Baseline': 0.8976,
    'WideNet Enhanced': 0.8859
}

# --- Instantiate models for param counts ---

```

```

models = {
    'ShallowNet Baseline': ShallowNet(input_size, 64, output_size, enhanced=False),
    'ShallowNet Enhanced': ShallowNet(input_size, 64, output_size, enhanced=True),
    'DeepNet Baseline': DeepNet(input_size, 128, output_size, enhanced=False),
    'DeepNet Enhanced': DeepNet(input_size, 128, output_size, enhanced=True),
    'WideNet Baseline': WideNet(input_size, 512, output_size, enhanced=False),
    'WideNet Enhanced': WideNet(input_size, 512, output_size, enhanced=True)
}

param_counts = {name: sum(p.numel() for p in model.parameters()) for name, model
in models.items()}

# --- Build DataFrame ---
df = pd.DataFrame({
    'Model': list(mse_values.keys()),
    'MSE': list(mse_values.values()),
    'Params': [param_counts[name] for name in mse_values.keys()]
})

# --- Sort by MSE ---
df_sorted = df.sort_values('MSE', ascending=True).reset_index(drop=True)

# --- Plot heatmaps ---
fig, axes = plt.subplots(1, 2, figsize=(12, 3))

# MSE heatmap
sns.heatmap(
    df_sorted[['MSE']].T,
    annot=df_sorted['MSE'].apply(lambda x: f"{x:.4f}").values.reshape(1, -1),
    fmt='',
    cmap='RdYlGn_r', # green=low MSE, red=high MSE
    cbar_kws={'label': 'MSE'},
    xticklabels=df_sorted['Model'],
    ax=axes[0]
)
axes[0].set_title("MSE (lower is better)")
axes[0].set_yticks([])

# Params heatmap
sns.heatmap(
    df_sorted[['Params']].T,
    annot=df_sorted['Params'].apply(lambda x: f"{x:,}").values.reshape(1, -1),
    fmt='',
    cmap='YlGnBu', # blue scale for compute cost
    cbar_kws={'label': 'Parameter Count'},
    xticklabels=df_sorted['Model'],
    ax=axes[1]
)
axes[1].set_title("Compute Cost (Params)")
axes[1].set_yticks([])

plt.tight_layout()

```

```
plt.show()
```

▼ THAT PROVES IT

Sure, DeepNet is fractionally better.. and 150-something more computing cost.

ShallowNet is a teeny bit slower.. yet super low in cost. So it's pretty obvious what the most efficient method should be here.

ShallowNet wins it!
