# Abstract

**Lucas Little**
**CSCA5622**
**University of Colorado Boulder**

This project, "NBA Playoff Predictor," uses machine learning to predict which NBA teams will make the playoffs. Data was collected from several sources, including:

- NBA Shots Dataset
- NBA Injury Stats Dataset
- NBA/ABA/BAA Team Statistics Dataset

The data includes information about teams, players, and injuries. After cleaning and organizing the data, we created useful features like team performance and player experience.

We tested different machine learning models, such as Logistic Regression, Random Forest, and XGBoost, to see which one worked best. Logistic Regression was the top performer, with 83.3% accuracy on the test data.

This project shows how data and machine learning can help predict playoff success. The final model is ready for use and can be improved further in the future.

# NBA Playoff Predictor - Data Collection

This notebook gathers historical NBA team stats to create a dataset for predicting playoff outcomes. The data comes from multiple Kaggle sources.

## Notebook Overview

1. **Setup and Imports**: Load libraries and prepare the environment.
2. **Kaggle Data Collection**: Download team stats from Kaggle.
3. **Collection Summary**: Review the collected data.

## Setup and Imports

This section sets up the environment by importing necessary libraries and enabling logging to monitor progress and troubleshoot issues during data collection.

```python
In [6]: import os
        import sys
        from datetime import datetime

        # Add the src directory to the Python path to access utility functions
        sys.path.append('..')

        # Import utility functions and classes for logging and progress tracking
        from src.data.utils import setup_logging, DataCollectionProgress
        from src.data.collectors.kaggle_collector import KaggleCollector

        # Set up logging to track the progress of the data collection process
        logger = setup_logging()

        # Initialize a progress tracker to monitor the status of each data collectic
        progress = DataCollectionProgress()

        # Initialize the KaggleCollector to manage dataset downloads
        kaggle = KaggleCollector('../data/raw/kaggle')
```

# Kaggle Data Collection

We use the Kaggle API to download historical team stats, handling each dataset individually for better control and visibility.

## NBA Shot Locations Dataset

This dataset provides detailed shot locations and types to analyze team shooting profiles.

```python
In [7]: # Download NBA Shots Dataset
        progress.add_task('download_shots', total_steps=1)
        progress.start_task('download_shots')

        try:
            result = kaggle.download_dataset('nba_shots', 'mexwell/nba-shots')
            if result['status'] == 'success':
                logger.info("Successfully downloaded NBA shots dataset")
                progress.complete_task('download_shots')
            else:
                logger.error(f"Failed to download NBA shots dataset: {result['error'
                progress.complete_task('download_shots', success=False)
        except Exception as e:
            logger.error(f"Error downloading NBA shots dataset: {str(e)}")
            progress.complete_task('download_shots', success=False, error=str(e))
```

```
2024-12-10 23:38:34 - INFO - Downloading dataset: nba_shots
Dataset URL: https://www.kaggle.com/datasets/mexwell/nba-shots
2024-12-10 23:38:40 - INFO - Successfully downloaded nba_shots
2024-12-10 23:38:40 - INFO - Successfully downloaded NBA shots dataset
```

# NBA Injury Statistics Dataset

This dataset includes injury data from 1951-2023, offering insights into how injuries impact team performance.

```python
In [8]:  # Download NBA Injury Stats Dataset
         progress.add_task('download_injuries', total_steps=1)
         progress.start_task('download_injuries')

         try:
             result = kaggle.download_dataset('nba_injuries', 'loganlauton/nba-injury
             if result['status'] == 'success':
                 logger.info("Successfully downloaded NBA injury stats dataset")
                 progress.complete_task('download_injuries')
             else:
                 logger.error(f"Failed to download NBA injury stats dataset: {result[
                 progress.complete_task('download_injuries', success=False)
         except Exception as e:
             logger.error(f"Error downloading NBA injury stats dataset: {str(e)}")
             progress.complete_task('download_injuries', success=False, error=str(e))
```

```
2024-12-10 23:38:40 - INFO - Downloading dataset: nba_injuries
Dataset URL: https://www.kaggle.com/datasets/loganlauton/nba-injury-stats-19
51-2023
2024-12-10 23:38:40 - INFO - Successfully downloaded nba_injuries
2024-12-10 23:38:40 - INFO - Successfully downloaded NBA injury stats datase
t
```

## NBA/ABA/BAA Team Statistics Dataset

This main dataset contains historical team stats from 1950 to today, including regular season data and advanced metrics.

```python
In [9]:  # Download NBA Team Stats Dataset
         progress.add_task('download_team_stats', total_steps=1)
         progress.start_task('download_team_stats')

         try:
             result = kaggle.download_dataset('nba_team_stats', 'sumitrodatta/nba-aba
             if result['status'] == 'success':
                 logger.info("Successfully downloaded NBA team stats dataset")
                 progress.complete_task('download_team_stats')
             else:
                 logger.error(f"Failed to download NBA team stats dataset: {result['e
                 progress.complete_task('download_team_stats', success=False)
         except Exception as e:
             logger.error(f"Error downloading NBA team stats dataset: {str(e)}")
             progress.complete_task('download_team_stats', success=False, error=str(e
```

```
2024-12-10 23:38:40 - INFO - Downloading dataset: nba_team_stats
Dataset URL: https://www.kaggle.com/datasets/sumitrodatta/nba-aba-baa-stats
```

## Collection Summary

This section reviews the data collection results, summarizing the total tasks completed and their status.

```python
In [10]:  # Get a summary of the data collection process
          summary = progress.get_summary()

          print(f"Data Collection Summary:")
          print(f"Total Tasks: {summary['total_tasks']}")
          print(f"Completed Successfully: {summary['completed_tasks']}")
          print(f"Failed: {summary['failed_tasks']}")
          print(f"Total Duration: {summary['duration']}")

          print("\nTask Details:")
          for name, task in summary['tasks'].items():
              status = task['status']
              duration = task['end_time'] - task['start_time'] if task['end_time'] and
              print(f"\n{name}:")
              print(f"  Status: {status}")
              print(f"  Duration: {duration}")
              if task['error']:
                  print(f"  Error: {task['error']}")
```

```
Data Collection Summary:
Total Tasks: 3
Completed Successfully: 3
Failed: 0
Total Duration: 0:00:07.317472

Task Details:

download_shots:
  Status: completed
  Duration: 0:00:06.069880

download_injuries:
  Status: completed
  Duration: 0:00:00.341802

download_team_stats:
  Status: completed
  Duration: 0:00:00.855184
```

# NBA Playoff Predictor - Data Cleaning

This notebook cleans and preprocesses NBA data from Kaggle to prepare it for feature engineering. It ensures standardized formats, handles missing values, and maintains

data quality across sources.

# Data Sources and Cleaning Goals

1. **NBA/ABA/BAA Stats (sumitrodatta)**

   - **Player Season Info**: Individual player stats per season
     - Standardize team names
     - Filter for NBA-only data
     - Add conference information
     - Handle missing values
   - **Team Stats Per Game**: Team performance metrics
     - Normalize team names
     - Add conference mappings
     - Ensure consistent stats

2. **NBA Injury Stats (loganlauton)**

   - Historical injury data (1951-2023)
     - Remove data before 2024
     - Normalize team names
     - Remove missing values
     - Add conference mappings
     - Create yearly injury summaries by team

In [16]:
```python
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import sys
from pathlib import Path

sys.path.append('..')

from src.data.cleaners.nba_data_cleaner import NBACleaner
from src.data.utils import setup_logging

logger = setup_logging()

sns.set_theme()

cleaner = NBACleaner()
```

# Clean NBA/ABA/BAA Stats

Process data from sumitrodatta's dataset by standardizing formats, adding conference details, and removing anomalies to ensure accurate analysis.

# Player Season Data Cleaning Script

This script cleans and standardizes NBA player season data for analysis. It performs the following steps:

1. **Load Data**: Reads the raw player season data from a CSV file.
2. **Filter by Season**: Removes records from seasons before 2004 to focus on recent data.
3. **Drop Unnecessary Columns**: Removes irrelevant columns like `player_id`, `seas_id`, `birth_year`, `pos`, and `lg`.
4. **Standardize Names**:
   - Standardizes player names using `NBACleaner`
   - Converts team names to uppercase and standardizes them
5. **Add Conference Information**:
   - Maps teams to their respective conferences
   - No historical alignment handling since we only use data from 2004 onward
6. **Convert Columns**:
   - Converts the `age` column to an integer type
   - Standardizes other numeric columns as needed
7. **Validate Data**:
   - Checks for missing values
   - Verifies conference mappings

## Output

- A cleaned dataset saved as `../data/processed/player_season.csv`
- Includes player statistics with team and conference information

In [17]:
```python
# Load and clean player season info data
logger.info("Loading player season info data...")
ps_df_raw = pd.read_csv('../data/raw/kaggle/sumitrodatta/nba-aba-baa-stats/F
logger.info(f"Initial player season info records: {len(ps_df_raw):,}")

# Make a copy of the raw data for processing
ps_df_processed = ps_df_raw.copy()

# Filter out data before 2004
logger.info("Filtering out records before 2004...")
ps_df_processed = ps_df_processed[ps_df_processed['season'] >= 2004]
logger.info(f"Records after filtering pre-2004 data: {len(ps_df_processed):,

# Drop unnecessary columns
columns_to_drop = ['player_id', 'seas_id', 'birth_year', 'pos', 'lg']
logger.info(f"Dropping unnecessary columns: {columns_to_drop}...")
ps_df_processed.drop(columns=columns_to_drop, inplace=True)
logger.info("Unnecessary columns dropped.")
```

```python
# Standardize player names
logger.info("Standardizing player names...")
ps_df_processed = cleaner.standardize_player_names(ps_df_processed)
logger.info("Player name standardization complete.")

# Standardize team names
logger.info("Standardizing team names: converting to uppercase and stripping
ps_df_processed['tm'] = ps_df_processed['tm'].str.strip().str.upper()
ps_df_processed = cleaner.standardize_team_names(ps_df_processed, ['tm'])
logger.info("Team name standardization complete.")

# Rename 'tm' column to 'team'
logger.info("Renaming 'tm' column to 'team'...")
ps_df_processed.rename(columns={'tm': 'team'}, inplace=True)
logger.info("Column 'tm' successfully renamed to 'team'.")

logger.info("Dropping rows where team equals 'TOT'...")
initial_rows = len(ps_df_processed)
ps_df_processed = ps_df_processed[ps_df_processed['team'] != 'TOT']
rows_dropped = initial_rows - len(ps_df_processed)
logger.info(f"Dropped {rows_dropped:,} rows with team 'TOT'")
logger.info(f"Remaining rows: {len(ps_df_processed):,}")

# Convert 'age' column to integer type
logger.info("Converting 'age' column to integer type...")
ps_df_processed['age'] = pd.to_numeric(ps_df_processed['age'], errors='coerc
logger.info("'age' column conversion to integer complete.")

# Add conference mappings
ps_df_processed = cleaner.add_conference_mappings(ps_df_processed, name_col=

unknown_teams = ps_df_processed[ps_df_processed['conference'] == 'Unknown']|
if len(unknown_teams) > 0:
    logger.warning(f"Found teams with unknown conference: {unknown_teams}")
else:
    logger.info("All teams successfully mapped to conferences")

logger.info("Dropping rows with unknown conferences...")
initial_rows = len(ps_df_processed)
ps_df_processed = ps_df_processed[ps_df_processed['conference'] != 'Unknown'
rows_dropped = initial_rows - len(ps_df_processed)
logger.info(f"Dropped {rows_dropped:,} rows with unknown conferences")
logger.info(f"Remaining rows: {len(ps_df_processed):,}")

# Check for NaN values
logger.info("Checking for NaN values in the dataset...")
nan_cols = ps_df_processed.columns[ps_df_processed.isna().any()].tolist()
if nan_cols:
    logger.warning(f"Found NaN values in the following columns: {nan_cols}")
    for col in nan_cols:
        nan_count = ps_df_processed[col].isna().sum()
        logger.warning(f"Column '{col}' has {nan_count:,} NaN values.")
else:
    logger.info("No NaN values found in the cleaned DataFrame.")
```

```python
# Display sample and save
print("\nSample of cleaned data:")
print(ps_df_processed.head())

output_path = '../data/processed/player_season.csv'
ps_df_processed.to_csv(output_path, index=False)
logger.info(f"Cleaned data saved to {output_path}")
```

```
2024-12-11 21:28:18 - INFO - Loading player season info data...
2024-12-11 21:28:18 - INFO - Initial player season info records: 32,358
2024-12-11 21:28:18 - INFO - Filtering out records before 2004...
2024-12-11 21:28:18 - INFO - Records after filtering pre-2004 data: 13,629
2024-12-11 21:28:18 - INFO - Dropping unnecessary columns: ['player_id', 'se
as_id', 'birth_year', 'pos', 'lg']...
2024-12-11 21:28:18 - INFO - Unnecessary columns dropped.
2024-12-11 21:28:18 - INFO - Standardizing player names...
2024-12-11 21:28:18 - INFO - Player name standardization complete.
2024-12-11 21:28:18 - INFO - Standardizing team names: converting to upperca
se and stripping whitespace...
2024-12-11 21:28:18 - INFO - Team name standardization complete.
2024-12-11 21:28:18 - INFO - Renaming 'tm' column to 'team'...
2024-12-11 21:28:18 - INFO - Column 'tm' successfully renamed to 'team'.
2024-12-11 21:28:18 - INFO - Dropping rows where team equals 'TOT'...
2024-12-11 21:28:18 - INFO - Dropped 1,348 rows with team 'TOT'
2024-12-11 21:28:18 - INFO - Remaining rows: 12,281
2024-12-11 21:28:18 - INFO - Converting 'age' column to integer type...
2024-12-11 21:28:18 - INFO - 'age' column conversion to integer complete.
2024-12-11 21:28:18 - INFO - All teams successfully mapped to conferences
2024-12-11 21:28:18 - INFO - Dropping rows with unknown conferences...
2024-12-11 21:28:18 - INFO - Dropped 0 rows with unknown conferences
2024-12-11 21:28:18 - INFO - Remaining rows: 12,281
2024-12-11 21:28:18 - INFO - Checking for NaN values in the dataset...
2024-12-11 21:28:18 - INFO - No NaN values found in the cleaned DataFrame.
2024-12-11 21:28:18 - INFO - Cleaned data saved to ../data/processed/player_
season.csv
```
```
Sample of cleaned data:
       season          player  age team  experience conference
18729    2004      Aaron McKie   31  PHI          10       EAST
18730    2004   Aaron Williams   32  BKN          10       EAST
18731    2004     Adonal Foyle   28  GSW           7       WEST
18732    2004    Adrian Griffin   29  HOU           5       WEST
18733    2004    Al Harrington   23  IND           6       EAST
```

This section visualizes the data cleaning process for the "Player Season Data" dataset, comparing the raw and processed data counts.

- **Dataset Name**: Player Season Data
- **Raw Data Count**: `raw_count`
- **Processed Data Count**: `processed_count`

In [18]:
```python
# Data for Player Season Data
dataset_name = "Player Season Data"
raw_count = len(ps_df_raw)
processed_count = len(ps_df_processed)
```

```python
# Data for plotting
stages = ['Raw Data (Combined)', 'Processed Data']
counts = [raw_count, processed_count]
colors = ['steelblue', 'orange']

# Create figure and axis
fig, ax = plt.subplots(figsize=(10, 4))

# Create horizontal bars
bars = ax.barh(stages, counts, color=colors)

# Add value labels
for bar in bars:
    width = bar.get_width()
    ax.text(width, bar.get_y() + bar.get_height() / 2,
            f'{int(width):,}',
            ha='left', va='center', fontweight='bold')

# Customize the chart
ax.set_title(f'Data Cleaning Funnel: {dataset_name}', pad=20)
ax.set_xlabel('Number of Records')

# Remove spines
ax.spines['top'].set_visible(False)
ax.spines['right'].set_visible(False)

# Add padding
ax.set_xlim(0, max(counts) * 1.15)

plt.tight_layout()
plt.show()
```
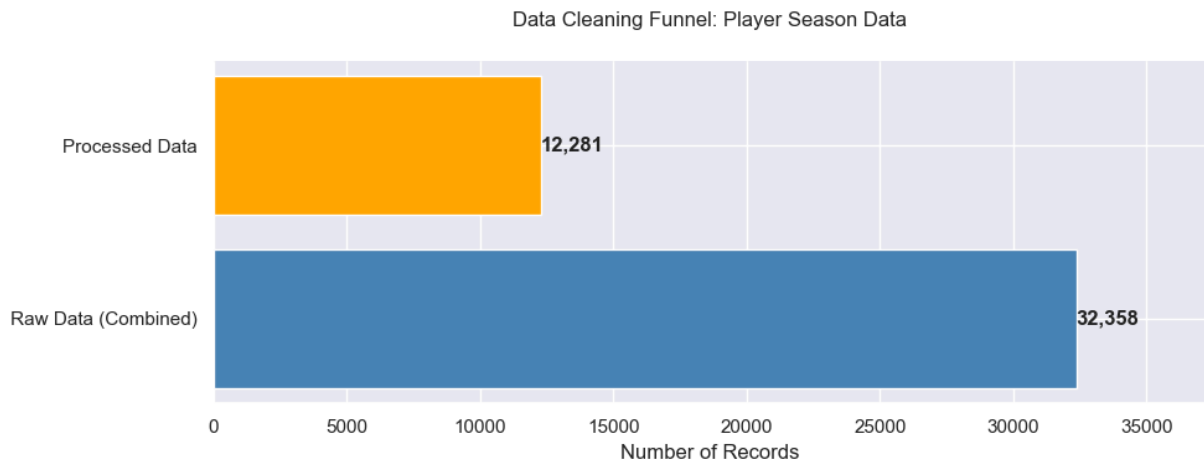


Data Cleaning Funnel: Player Season Data

# Team Statistics Data Cleaning Script

This script processes and cleans NBA team statistics data for analysis. It performs the following steps:

1. **Load Data**: Reads the raw team statistics data.

2. **Filter by Season**: Removes records prior to 2004.
3. **Remove League Averages**: Excludes league average entries.
4. **Standardize Team Names**:
   - Converts team names to uppercase
   - Maps team names to standardized codes
5. **Add Conference Information**:
   - Maps teams to their respective conferences
   - Validates conference assignments
6. **Convert Data Types**:
   - Handles numeric columns
   - Converts percentages to decimals
7. **Validate Data**:
   - Checks for missing values
   - Verifies data consistency

## Output

- A cleaned dataset saved as `../data/processed/team_stats.csv`
- Includes team statistics with conference information

In [19]:
```python
# Load and clean team stats data
logger.info("Loading team stats data...")
ts_df_raw = pd.read_csv('../data/raw/kaggle/sumitrodatta/nba-aba-baa-stats/T
logger.info(f"Initial team stats records: {len(ts_df_raw):,}")

# Make a copy of the raw data for processing
ts_df_processed = ts_df_raw.copy()

# Filter out data before 2004
logger.info("Filtering out records before 2004...")
ts_df_processed = ts_df_processed[ts_df_processed['season'] >= 2004]
logger.info(f"Records after filtering pre-2004 data: {len(ts_df_processed):,

# Remove League Average entries
logger.info("Removing 'League Average' entries from the data...")
ts_df_processed = ts_df_processed[~ts_df_processed['team'].str.contains('Lea
logger.info(f"Records after removing 'League Average' entries: {len(ts_df_pr

# Standardize team names
logger.info("Standardizing team names...")
ts_df_processed['team'] = ts_df_processed['team'].str.strip().str.upper()
ts_df_processed = cleaner.standardize_team_names(ts_df_processed, ['team'])
logger.info("Team name standardization complete.")

# Convert percentage strings to decimals
logger.info("Converting percentage strings to decimal values...")
ts_df_processed = cleaner.convert_percentages(ts_df_processed)
logger.info("Percentage conversion complete.")

# Add conference mappings
```

```python
ts_df_processed = cleaner.add_conference_mappings(ts_df_processed, name_col=

unknown_teams = ts_df_processed[ts_df_processed['conference'] == 'Unknown'][
if len(unknown_teams) > 0:
    logger.warning(f"Found teams with unknown conference: {unknown_teams}")
else:
    logger.info("All teams successfully mapped to conferences")

# Check for NaN values
logger.info("Checking for NaN values in the dataset...")
nan_cols = ts_df_processed.columns[ts_df_processed.isna().any()].tolist()
if nan_cols:
    logger.warning(f"Found NaN values in the following columns: {nan_cols}")
    for col in nan_cols:
        nan_count = ts_df_processed[col].isna().sum()
        logger.warning(f"Column '{col}' has {nan_count:,} NaN values.")
else:
    logger.info("No NaN values found in the cleaned DataFrame.")

# Display sample and save
print("\nSample of cleaned data:")
print(ts_df_processed.head())

output_path = '../data/processed/team_stats.csv'
ts_df_processed.to_csv(output_path, index=False)
logger.info(f"Cleaned data saved to {output_path}")
```

```
2024-12-11 21:28:18 - INFO - Loading team stats data...
2024-12-11 21:28:18 - INFO - Initial team stats records: 1,876
2024-12-11 21:28:18 - INFO - Filtering out records before 2004...
2024-12-11 21:28:18 - INFO - Records after filtering pre-2004 data: 681
2024-12-11 21:28:18 - INFO - Removing 'League Average' entries from the dat
a...
2024-12-11 21:28:18 - INFO - Records after removing 'League Average' entrie
s: 659
2024-12-11 21:28:18 - INFO - Standardizing team names...
2024-12-11 21:28:18 - INFO - Team name standardization complete.
2024-12-11 21:28:18 - INFO - Converting percentage strings to decimal value
s...
2024-12-11 21:28:18 - INFO - Percentage conversion complete.
2024-12-11 21:28:18 - INFO - All teams successfully mapped to conferences
2024-12-11 21:28:18 - INFO - Checking for NaN values in the dataset...
2024-12-11 21:28:18 - INFO - No NaN values found in the cleaned DataFrame.
2024-12-11 21:28:18 - INFO - Cleaned data saved to ../data/processed/team_st
ats.csv
```

```
Sample of cleaned data:
   season   lg team abbreviation  playoffs     g  mp_per_game  fg_per_game
\
0    2025  NBA  ATL          ATL     False  21.0        240.0         42.4
1    2025  NBA  BOS          BOS     False  19.0        243.9         41.8
2    2025  NBA  BKN          BRK     False  20.0        242.5         39.4
3    2025  NBA  CHI          CHI     False  21.0        240.0         43.1
4    2025  NBA  CHA          CHO     False  20.0        242.5         38.6

   fga_per_game  fg_percent  ...  orb_per_game  drb_per_game  trb_per_game
\
0          91.5       0.463  ...          12.3          32.9          45.2
1          90.3       0.464  ...          10.1          33.2          43.3
2          84.0       0.468  ...           8.5          29.8          38.3
3          90.8       0.475  ...           8.9          35.2          44.1
4          91.0       0.424  ...          13.4          32.4          45.8

   ast_per_game  stl_per_game  blk_per_game  tov_per_game  pf_per_game  \
0          29.8           9.9           5.6          16.4         18.8
1          25.9           7.3           5.3          11.6         16.6
2          26.4           6.4           3.7          14.6         22.3
3          28.7           7.0           4.6          15.4         18.4
4          22.9           8.1           5.2          16.2         21.3

   pts_per_game  conference
0         116.1        EAST
1         121.2        EAST
2         111.8        EAST
3         118.5        EAST
4         107.5        EAST

[5 rows x 29 columns]
```

This section visualizes the data cleaning process for the "Team Statistics Data" dataset, comparing the raw and processed data counts.

- **Dataset Name**: Team Statistics Data
- **Raw Data Count**: `raw_count`
- **Processed Data Count**: `processed_count`

In [20]:
```python
# Data for Team Statistics data
dataset_name = "Team Statistics data"
raw_count = len(ts_df_raw)
processed_count = len(ts_df_processed)

# Data for plotting
stages = ['Raw Data (Combined)', 'Processed Data']
counts = [raw_count, processed_count]
colors = ['steelblue', 'orange']

# Create figure and axis
fig, ax = plt.subplots(figsize=(10, 4))

# Create horizontal bars
bars = ax.barh(stages, counts, color=colors)
```

```python
# Add value labels
for bar in bars:
    width = bar.get_width()
    ax.text(width, bar.get_y() + bar.get_height() / 2,
            f'{int(width):,}',
            ha='left', va='center', fontweight='bold')

# Customize the chart
ax.set_title(f'Data Cleaning Funnel: {dataset_name}', pad=20)
ax.set_xlabel('Number of Records')

# Remove spines
ax.spines['top'].set_visible(False)
ax.spines['right'].set_visible(False)

# Add padding
ax.set_xlim(0, max(counts) * 1.15)

plt.tight_layout()
plt.show()
```
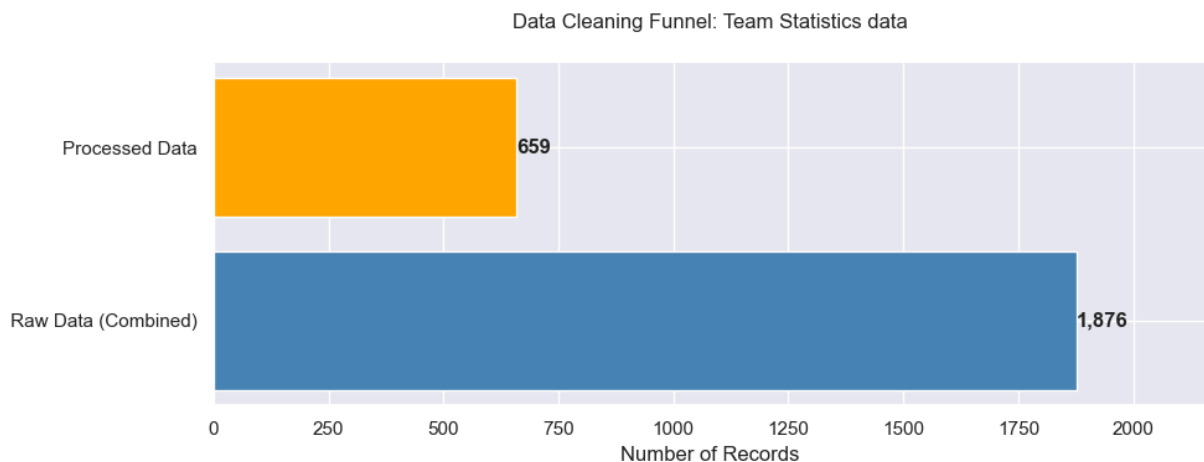


Data Cleaning Funnel: Team Statistics data

# Injury Data Cleaning and Summary Script

This script processes NBA player injury data to generate team-level injury summaries. It performs the following steps:

1. **Load Data**: Reads raw injury data.
2. **Clean Dates**: Converts date columns to a consistent datetime format.
3. **Filter Data**: Removes records before 2004.
4. **Standardize Teams**:
   - Converts team names to a standardized format
   - Maps historical team names to current ones
5. **Add Conference Information**:
   - Maps teams to their respective conferences
   - Validates conference assignments

6. **Create Summary**:
   - Groups data by year and team
   - Counts injuries per team-season

## Outputs

- **Injury Summary**: A CSV file saved as
  `../data/processed/injuries_summary.csv`
  - Includes year, team, conference, and injury count

In [21]:
```python
# Load and clean injury data
logger.info("Loading player injury data...")
injury_df_raw = pd.read_csv('../data/raw/kaggle/loganlauton/nba-injury-stats
logger.info(f"Initial player injury records: {len(injury_df_raw):,}")

# Make a copy of the raw data for processing
injury_df_processed = injury_df_raw.copy()

# Convert dates to datetime
logger.info("Converting 'Date' column to datetime format...")
injury_df_processed = cleaner.handle_dates(injury_df_processed, ['Date'])
logger.info("Date conversion complete.")

# Filter out data before 2004
logger.info("Filtering out records before 2004...")
injury_df_processed = injury_df_processed[injury_df_processed['Date'] >= '20
logger.info(f"Records after filtering pre-2004 data: {len(injury_df_processe

# Drop unnecessary columns
columns_to_drop = ['Unnamed: 0', 'Acquired', 'Relinquished', 'Notes']
logger.info(f"Dropping unnecessary columns: {columns_to_drop}...")
injury_df_processed.drop(columns=columns_to_drop, inplace=True)
logger.info("Unnecessary columns dropped.")

# Standardize team names
logger.info("Standardizing team names...")
injury_df_processed['Team'] = injury_df_processed['Team'].str.strip().str.up
injury_df_processed = cleaner.standardize_team_names(injury_df_processed, ['
logger.info("Team name standardization complete.")

# Extract year and create summary
logger.info("Creating injury summary...")
injury_df_processed['Year'] = injury_df_processed['Date'].dt.year
injury_summary_df = injury_df_processed.groupby(['Year', 'Team']).size().res

# Add conference mappings
injury_summary_df = cleaner.add_conference_mappings(injury_summary_df, name_

unknown_teams = injury_summary_df[injury_summary_df['conference'] == 'Unknow
if len(unknown_teams) > 0:
    logger.warning(f"Found teams with unknown conference: {unknown_teams}")
else:
    logger.info("All teams successfully mapped to conferences")
```

```python
# Standardize column names
logger.info("Converting column names to lowercase...")
injury_summary_df.columns = injury_summary_df.columns.str.lower()

# Display sample and save
print("\nSample of injury summary:")
print(injury_summary_df.head())

output_path = '../data/processed/injuries_summary.csv'
injury_summary_df.to_csv(output_path, index=False)
logger.info(f"Injury summary saved to {output_path}")
```

```
2024-12-11 21:28:18 - INFO - Loading player injury data...
2024-12-11 21:28:19 - INFO - Initial player injury records: 37,667
2024-12-11 21:28:19 - INFO - Converting 'Date' column to datetime format...
2024-12-11 21:28:19 - INFO - Date conversion complete.
2024-12-11 21:28:19 - INFO - Filtering out records before 2004...
2024-12-11 21:28:19 - INFO - Records after filtering pre-2004 data: 31,385
2024-12-11 21:28:19 - INFO - Dropping unnecessary columns: ['Unnamed: 0', 'A
cquired', 'Relinquished', 'Notes']...
2024-12-11 21:28:19 - INFO - Unnecessary columns dropped.
2024-12-11 21:28:19 - INFO - Standardizing team names...
2024-12-11 21:28:19 - INFO - Team name standardization complete.
2024-12-11 21:28:19 - INFO - Creating injury summary...
2024-12-11 21:28:19 - INFO - All teams successfully mapped to conferences
2024-12-11 21:28:19 - INFO - Converting column names to lowercase...
2024-12-11 21:28:19 - INFO - Injury summary saved to ../data/processed/injur
ies_summary.csv
```
```
Sample of injury summary:
   year team  count conference
0  2004  ATL     20       EAST
1  2004  BKN     21       EAST
2  2004  BOS     16       EAST
3  2004  CHA     36       EAST
4  2004  CHI     25       EAST
```

# Data for Player Injury Data

This section visualizes the data cleaning process for the "Player Injury Data" dataset, comparing the raw and processed data counts.

- **Dataset Name**: Player Injury Data
- **Raw Data Count**: `raw_count`
- **Processed Data Count**: `processed_count`

In [22]:
```python
# Data for Player Injury Data
dataset_name = "Player Injury Data"
raw_count = len(injury_df_raw)
processed_count = len(injury_df_processed)

# Data for plotting
stages = ['Raw Data (Combined)', 'Processed Data']
```

```
counts = [raw_count, processed_count]
colors = ['steelblue', 'orange']

# Create figure and axis
fig, ax = plt.subplots(figsize=(10, 4))

# Create horizontal bars
bars = ax.barh(stages, counts, color=colors)

# Add value labels
for bar in bars:
    width = bar.get_width()
    ax.text(width, bar.get_y() + bar.get_height() / 2,
            f'{int(width):,}',
            ha='left', va='center', fontweight='bold')

# Customize the chart
ax.set_title(f'Data Cleaning Funnel: {dataset_name}', pad=20)
ax.set_xlabel('Number of Records')

# Remove spines
ax.spines['top'].set_visible(False)
ax.spines['right'].set_visible(False)

# Add padding
ax.set_xlim(0, max(counts) * 1.15)

plt.tight_layout()
plt.show()
```
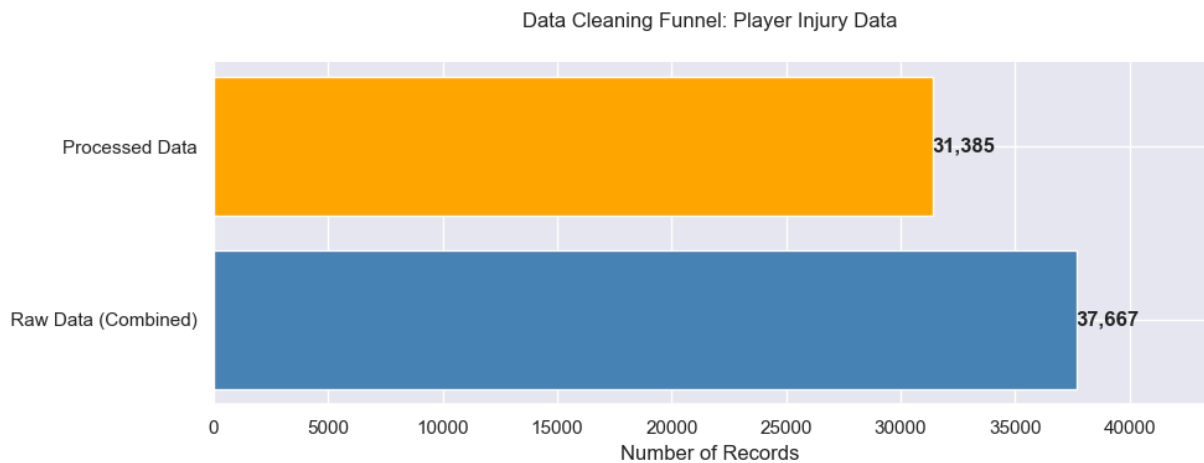


Data Cleaning Funnel: Player Injury Data

# Conclusions

This section summarizes the key findings and outcomes of the data cleaning process.

- **Data Cleaning Success**: The raw data from multiple sources has been successfully cleaned, standardized, and processed, ensuring consistency and accuracy for further analysis.

- **Data Quality Improvement**: Missing values have been addressed, team names standardized, and conference mappings added to all datasets.
- **Summary of Key Metrics**:
  - The number of records in the raw data has been reduced as we focused on relevant data (post-2004) and removed unnecessary columns.
  - Data from multiple sources, including player season statistics, team stats, and injury data, is now consistent and ready for feature engineering and analysis.
- **Next Steps**:
  - The cleaned datasets will be used for feature engineering and training the NBA Playoff Predictor model.
  - Further analysis will focus on identifying key predictors for playoff success, using the processed player and team statistics, along with injury data summaries.

# NBA Playoff Predictor - Feature Engineering

This notebook creates features for predicting NBA playoff outcomes using historical data. It generates factors that represent different parts of team performance, player makeup, and injury patterns that could affect a team's chances of making the playoffs.

## Data Sources

1. **Team Statistics** (`team_stats.csv`):

   - Includes team performance metrics per season, such as scoring, rebounding, assists, and defense.
   - Contains information about playoff qualification (target variable).
2. **Player Statistics** (`player_season.csv`):

   - Contains individual player data per season, including experience and age.
   - Helps analyze the team's composition.
3. **Injury Summary** (`injuries_summary.csv`):

   - Provides team-level injury counts by year.
   - Used to assess how injuries impact team performance.

```
In [27]:   # Data manipulation
           import pandas as pd
           import numpy as np

           # System utilities
           import sys
           from pathlib import Path
           import json
           from datetime import datetime
```

```python
# Visualization
import matplotlib.pyplot as plt
import seaborn as sns

# Add the src directory to the path
sys.path.append('..')
from src.features.feature_builder import FeatureBuilder
from src.data.utils import setup_logging

logger = setup_logging()

# Configure pandas display
pd.set_option('display.max_columns', None)
pd.set_option('display.max_rows', 100)
```

# Load and Explore Data

In this section, we load and preview the datasets used for our analysis. Each dataset offers valuable insights that will contribute to the feature engineering process:

## 1. Team Statistics (`team_stats.csv`)

This dataset includes season-level performance data for each team, such as:

- **Performance Metrics**: Points per game, field goal percentage, assists, and rebounds.
- **Playoff Qualification**: A binary target variable, `playoffs`, indicating if a team made the playoffs in a given season.

## 2. Player Statistics (`player_season.csv`)

This dataset contains player-level data for each season, offering:

- **Demographics**: Player age and experience.
- **Team Composition**: Information about roster depth and the distribution of experience across players.

## 3. Injury Summary (`injuries_summary.csv`)

This dataset provides aggregated team-level injury counts by season, allowing:

- **Injury Impact Analysis**: Insights into how injuries affect team performance and playoff chances.
- **Roster Stability**: Evaluation of player availability across seasons.

Each dataset will be explored briefly to ensure we understand its structure and content.

```
In [28]:  # Load all data sources
          data_dir = '../data/processed'

          # Load Team Statistics
          print("Loading Team Statistics...")
          team_stats = pd.read_csv(f'{data_dir}/team_stats.csv')
          print(f"Team Stats Shape: {team_stats.shape}")
          display(team_stats[['season', 'team', 'playoffs', 'pts_per_game', 'fg_percen

          # Load Player Statistics
          print("\nLoading Player Statistics...")
          player_stats = pd.read_csv(f'{data_dir}/player_season.csv')
          print(f"Player Stats Shape: {player_stats.shape}")
          display(player_stats[['season', 'player', 'team', 'age', 'experience']].head

          # Load Injury Summary
          print("\nLoading Injury Summary...")
          injuries = pd.read_csv(f'{data_dir}/injuries_summary.csv')
          print(f"Injuries Shape: {injuries.shape}")
          display(injuries.head())

          # Summarize datasets
          print("\nTeam Statistics Summary:")
          print(team_stats.info())

          print("\nPlayer Statistics Summary:")
          print(player_stats.info())

          print("\nInjury Summary Summary:")
          print(injuries.info())
```

```
Loading Team Statistics...
Team Stats Shape: (659, 29)
```

|   | season | team | playoffs | pts_per_game | fg_percent |
|---|--------|------|----------|--------------|------------|
| **0** | 2025 | ATL | False | 116.1 | 0.463 |
| **1** | 2025 | BOS | False | 121.2 | 0.464 |
| **2** | 2025 | BKN | False | 111.8 | 0.468 |
| **3** | 2025 | CHI | False | 118.5 | 0.475 |
| **4** | 2025 | CHA | False | 107.5 | 0.424 |

```
Loading Player Statistics...
Player Stats Shape: (12281, 6)
```

|   | season | player | team | age | experience |
|---|--------|--------|------|-----|-----------|
| 0 | 2004 | Aaron McKie | PHI | 31 | 10 |
| 1 | 2004 | Aaron Williams | BKN | 32 | 10 |
| 2 | 2004 | Adonal Foyle | GSW | 28 | 7 |
| 3 | 2004 | Adrian Griffin | HOU | 29 | 5 |
| 4 | 2004 | Al Harrington | IND | 23 | 6 |

```
Loading Injury Summary...
Injuries Shape: (591, 4)
```

|   | year | team | count | conference |
|---|------|------|-------|-----------|
| 0 | 2004 | ATL | 20 | EAST |
| 1 | 2004 | BKN | 21 | EAST |
| 2 | 2004 | BOS | 16 | EAST |
| 3 | 2004 | CHA | 36 | EAST |
| 4 | 2004 | CHI | 25 | EAST |

```
Team Statistics Summary:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 659 entries, 0 to 658
Data columns (total 29 columns):
 #   Column          Non-Null Count  Dtype
---  ------          --------------  -----
 0   season           659 non-null   int64
 1   lg               659 non-null   object
 2   team             659 non-null   object
 3   abbreviation     659 non-null   object
 4   playoffs         659 non-null   bool
 5   g                659 non-null   float64
 6   mp_per_game      659 non-null   float64
 7   fg_per_game      659 non-null   float64
 8   fga_per_game     659 non-null   float64
 9   fg_percent       659 non-null   float64
 10  x3p_per_game     659 non-null   float64
 11  x3pa_per_game    659 non-null   float64
 12  x3p_percent      659 non-null   float64
 13  x2p_per_game     659 non-null   float64
 14  x2pa_per_game    659 non-null   float64
 15  x2p_percent      659 non-null   float64
 16  ft_per_game      659 non-null   float64
 17  fta_per_game     659 non-null   float64
 18  ft_percent       659 non-null   float64
 19  orb_per_game     659 non-null   float64
 20  drb_per_game     659 non-null   float64
 21  trb_per_game     659 non-null   float64
 22  ast_per_game     659 non-null   float64
 23  stl_per_game     659 non-null   float64
 24  blk_per_game     659 non-null   float64
 25  tov_per_game     659 non-null   float64
 26  pf_per_game      659 non-null   float64
 27  pts_per_game     659 non-null   float64
 28  conference       659 non-null   object
dtypes: bool(1), float64(23), int64(1), object(4)
memory usage: 144.9+ KB
None

Player Statistics Summary:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 12281 entries, 0 to 12280
Data columns (total 6 columns):
 #   Column      Non-Null Count  Dtype
---  ------      --------------  -----
 0   season      12281 non-null  int64
 1   player      12281 non-null  object
 2   age         12281 non-null  int64
 3   team        12281 non-null  object
 4   experience  12281 non-null  int64
 5   conference  12281 non-null  object
dtypes: int64(3), object(3)
memory usage: 575.8+ KB
None

Injury Summary Summary:
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 591 entries, 0 to 590
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype
---  ------      --------------  -----
 0   year        591 non-null    int64
 1   team        591 non-null    object
 2   count       591 non-null    int64
 3   conference  591 non-null    object
dtypes: int64(2), object(2)
memory usage: 18.6+ KB
None
```

# Create Team Features

To improve our predictions, we create team-level features using the `team_stats` dataset. These features capture important aspects of team performance, grouped into the following categories:

## Efficiency Metrics

These metrics evaluate a team's shooting and rebounding efficiency:

- **True Shooting Percentage (TS%)**: Measures overall scoring efficiency, accounting for free throws and three-pointers.
- **Effective Field Goal Percentage (eFG%)**: Adjusts for the extra value of three-point shots.
- **Offensive Rebound Percentage (OREB%)**: Indicates how well a team secures offensive rebounds.

## Ball Movement and Control

These metrics assess the quality of a team's ball movement and possession management:

- **Assist-to-Turnover Ratio (AST/TO)**: Shows the efficiency of passing relative to turnovers.
- **Assist Ratio (AST%)**: The percentage of field goals assisted by teammates.

## Defensive Impact

These metrics reflect defensive effectiveness:

- **Stocks per Game**: The total number of steals and blocks per game.
- **Defensive Rating**: A comprehensive measure of a team's defensive performance.

## Possession and Pace

Metrics related to the speed and efficiency of play:

- **Possessions per Game**: Estimates the team's pace of play.
- **Offensive Efficiency (OffEff)**: Points scored per possession.
- **Three Point Rate (3P%)**: Shows how often a team relies on three-point shots.
- **Free Throw Rate (FTR)**: Evaluates a team's ability to draw fouls and get to the free throw line.

## Sample and Summary

Below, we display a subset of the engineered features along with their summary statistics to better understand their distributions.

In [29]:
```python
# Create team features
print("Creating team performance features...")
team_features = FeatureBuilder().create_team_features(team_stats)

# Display progress and sample
print(f"Created {len(team_features.columns)} team features for {len(team_fea
print("\nSample of engineered team features:")
feature_sample = [
    'true_shooting_pct', 'efg_pct', 'ast_to_ratio',
    'stocks_per_game', 'off_efficiency', 'efficiency_rating'
]
display(team_features[['team', 'season'] + feature_sample].head())

# Display summary statistics for selected features
print("\nSummary statistics for key features:")
display(team_features[feature_sample].describe().transpose())

# Visualize feature distributions (optional)
team_features[feature_sample].hist(bins=15, figsize=(15, 10), edgecolor='bla
plt.suptitle("Distribution of Key Engineered Features", fontsize=16)
plt.show()
```

```
Creating team performance features...
Created 41 team features for 659 seasons.

Sample of engineered team features:
```
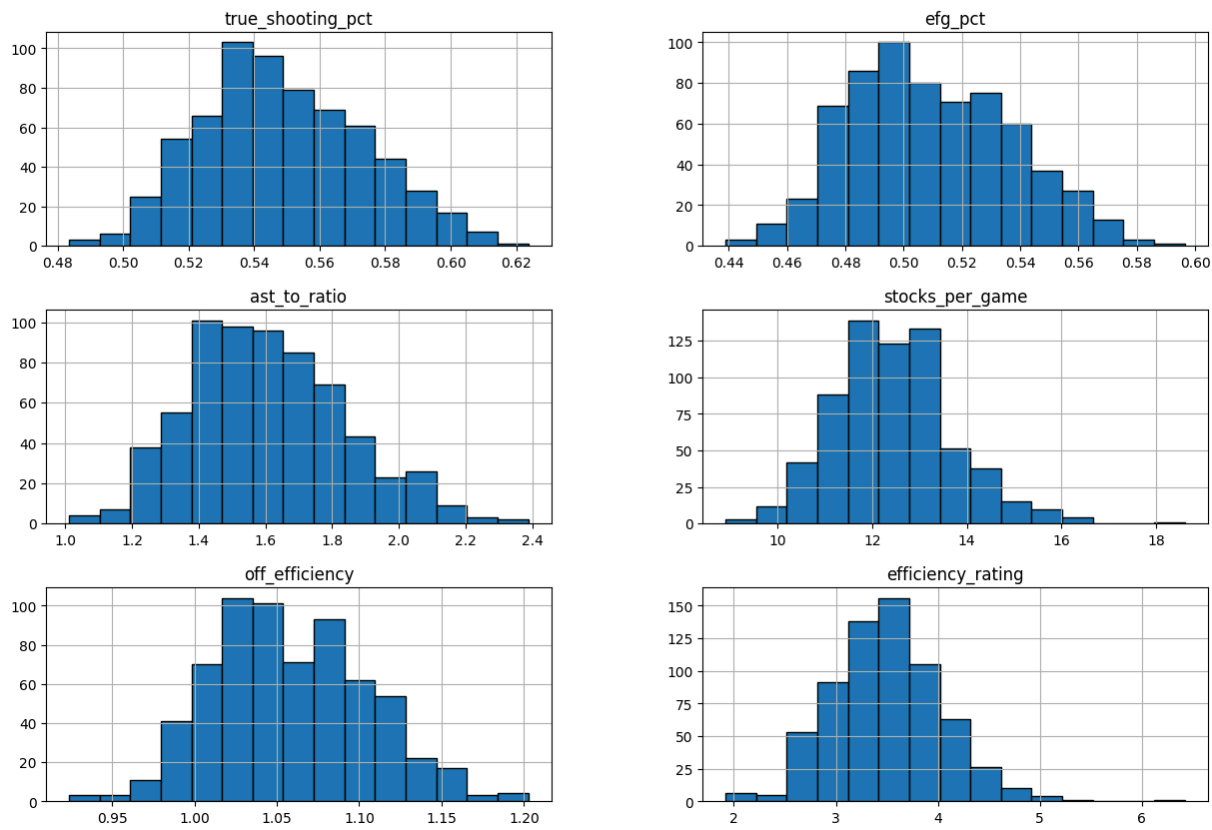
|   | team | season | true_shooting_pct | efg_pct | ast_to_ratio | stocks_per_game | off_effi |
|---|------|--------|-------------------|---------|--------------|-----------------|----------|
| 0 | ATL | 2025 | 0.568538 | 0.531694 | 1.817073 | 15.5 | 1.0 |
| 1 | BOS | 2025 | 0.604260 | 0.569214 | 2.232759 | 12.6 | 1.1 |
| 2 | BKN | 2025 | 0.598399 | 0.563095 | 1.808219 | 10.1 | 1.1 |
| 3 | CHI | 2025 | 0.594617 | 0.564978 | 1.863636 | 11.6 | 1.1 |
| 4 | CHA | 2025 | 0.543368 | 0.513187 | 1.413580 | 13.3 | 1.0 |

```
Summary statistics for key features:
```

|  | count | mean | std | min | 25% | 50% | 7 |
|---|---|---|---|---|---|---|---|
| **true_shooting_pct** | 659.0 | 0.548909 | 0.024707 | 0.483505 | 0.531726 | 0.546511 | 0.566 |
| **efg_pct** | 659.0 | 0.509504 | 0.027871 | 0.438903 | 0.488742 | 0.506208 | 0.530 |
| **ast_to_ratio** | 659.0 | 1.604520 | 0.231319 | 1.011299 | 1.435099 | 1.583333 | 1.755 |
| **stocks_per_game** | 659.0 | 12.445827 | 1.266134 | 8.900000 | 11.600000 | 12.400000 | 13.250 |
| **off_efficiency** | 659.0 | 1.058324 | 0.046578 | 0.923881 | 1.024216 | 1.053412 | 1.090 |
| **efficiency_rating** | 659.0 | 3.499821 | 0.534022 | 1.918881 | 3.142246 | 3.490173 | 3.841 |

Distribution of Key Engineered Features



# Create Player Features

This section focuses on creating features based on player statistics and injury data. These features provide insights into team composition and overall health, which are key factors in predicting a team's success.

## Experience and Age

We calculate metrics that capture the experience and age distribution of the team's roster:

- **Experience Metrics**: The average, maximum, and minimum years of experience for the team.
- **Age Demographics**: Includes the average, oldest, and youngest player ages.
- **Roster Size**: The total number of players on the team in a given season.

## Health Impact

Injury data is used to evaluate the impact of player availability:

- **Injury Count**: The total number of injuries affecting the team in a season.
- This metric helps assess the stability and depth of the roster.

## Sample and Summary

Below, we show a sample of the engineered features, followed by summary statistics and distributions. These metrics provide a clear picture of the team's composition and health.

In [30]:
```python
# Create player and injury features
print("Creating player composition and injury features...")
player_features = FeatureBuilder().create_player_features(player_stats, inju

# Display progress and sample
print(f"Created {len(player_features.columns)} player/injury features for {l
feature_sample = [
    'avg_experience', 'roster_size', 'avg_age',
    'max_age', 'min_age', 'count'
]

# Show sample of engineered features
print("\nSample of engineered player features:")
display(player_features[['team', 'season'] + feature_sample].head())

# Display summary statistics
print("\nSummary statistics for key features:")
display(player_features[feature_sample].describe().transpose())

# Visualize distributions of selected features
player_features[feature_sample].hist(bins=15, figsize=(15, 10), edgecolor='b
plt.suptitle("Distribution of Player Composition and Injury Features", fonts
plt.show()
```

Creating player composition and injury features...
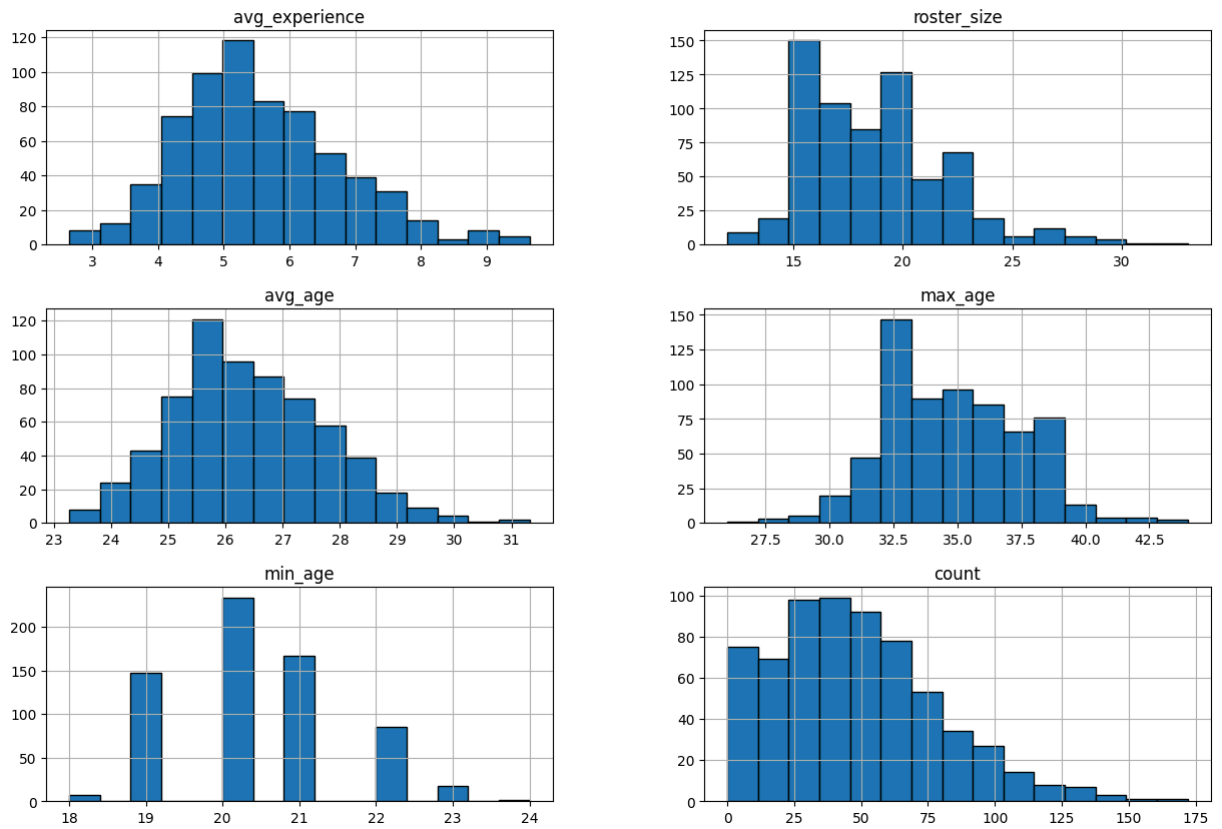Created 11 player/injury features for 659 seasons.

Sample of engineered player features:

|   | team | season | avg_experience | roster_size | avg_age | max_age | min_age | count |
|---|------|--------|----------------|-------------|---------|---------|---------|-------|
| 0 | ATL | 2004 | 5.043478 | 23 | 26.695652 | 32 | 21 | 20.0 |
| 1 | ATL | 2005 | 6.750000 | 20 | 28.000000 | 42 | 19 | 37.0 |
| 2 | ATL | 2006 | 3.533333 | 15 | 24.133333 | 32 | 19 | 32.0 |
| 3 | ATL | 2007 | 4.000000 | 19 | 24.631579 | 32 | 20 | 78.0 |
| 4 | ATL | 2008 | 4.875000 | 16 | 25.125000 | 33 | 21 | 39.0 |

Summary statistics for key features:

|  | count | mean | std | min | 25% | 50% | |
|---|-------|------|-----|-----|-----|-----|---|
| avg_experience | 659.0 | 5.565131 | 1.233115 | 2.652174 | 4.684211 | 5.388889 | 6.312 |
| roster_size | 659.0 | 18.635812 | 3.199685 | 12.000000 | 16.000000 | 18.000000 | 20.500 |
| avg_age | 659.0 | 26.403731 | 1.330941 | 23.263158 | 25.500000 | 26.230769 | 27.303 |
| max_age | 659.0 | 34.705615 | 2.648877 | 26.000000 | 33.000000 | 35.000000 | 36.500 |
| min_age | 659.0 | 20.356601 | 1.090996 | 18.000000 | 20.000000 | 20.000000 | 21.000 |
| count | 659.0 | 47.570561 | 31.148909 | 0.000000 | 25.000000 | 45.000000 | 65.000 |

Distribution of Player Composition and Injury Features



# Create Conference Features

To assess team performance within their respective conferences, we create several conference-specific features using the `team_stats` dataset. These features offer valuable context for comparing a team's performance against their conference peers and the overall standings.

## Conference Performance Metrics

The following metrics are calculated:

- **Points vs Conference Average**: Compares team scoring to the average points scored by other teams in the conference.
- **Wins vs Conference Average**: Compares a team's win total to the average wins in the conference.
- **Conference Rank**: Shows a team's position in the conference standings (1 = best).
- **Games Behind**: Indicates how far behind a team is from the conference leader in terms of wins.

```python
In [31]:  # Create conference features
print("Creating conference-based features...")
conference_features = FeatureBuilder().create_conference_features(team_stats

# Display progress and sample
print(f"Created {len(conference_features.columns)} conference features for {
feature_sample = [
    'conf_rank', 'pts_behind_leader', 'pts_behind_8th',
    'pts_vs_conf_avg', 'games_vs_conf_avg'
]

# Show sample of engineered features
print("\nSample of engineered conference features:")
display(conference_features[['team', 'season'] + feature_sample].head())

# Display summary statistics
print("\nSummary statistics for key features:")
display(conference_features[feature_sample].describe().transpose())

# Visualize distributions of selected features
conference_features[feature_sample].hist(bins=15, figsize=(15, 10), edgecolc
plt.suptitle("Distribution of Conference-Based Features", fontsize=16)
plt.show()
```

```
Creating conference-based features...
Created 8 conference features for 659 seasons.

Sample of engineered conference features:
```
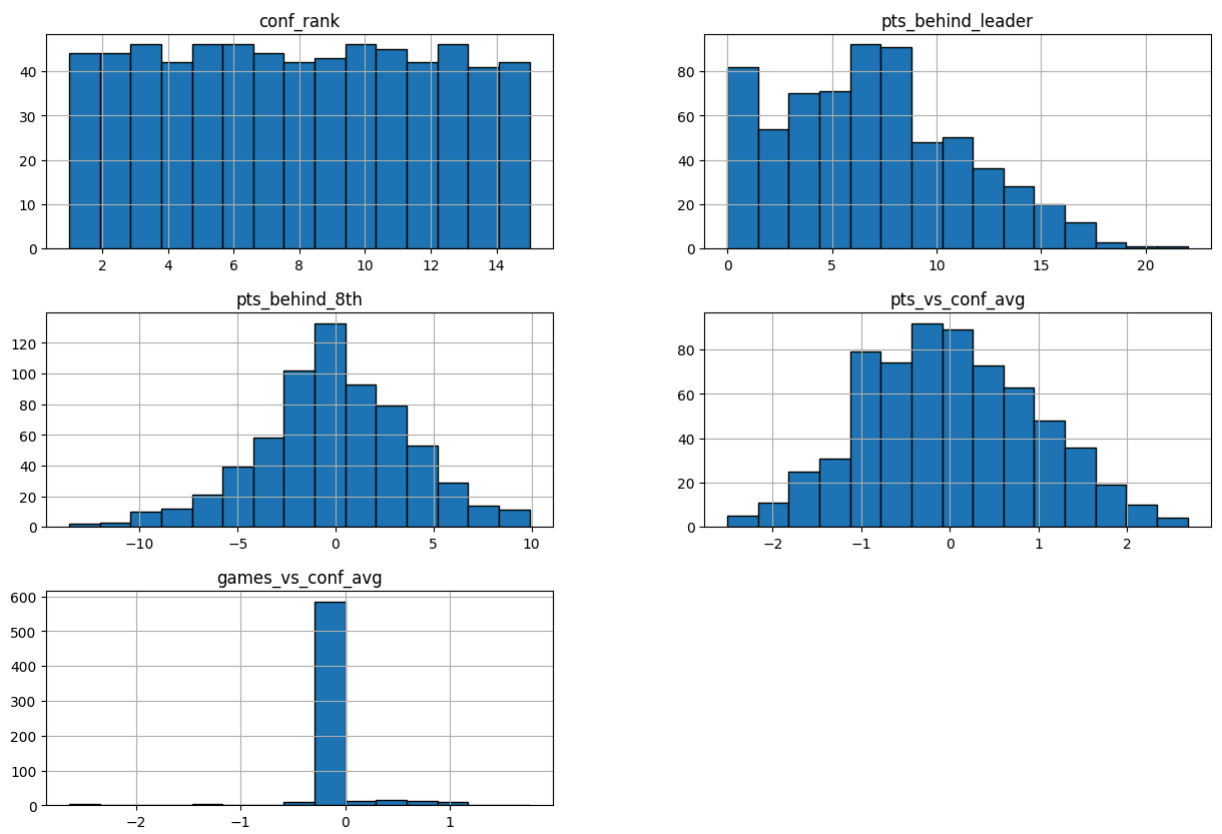
| | team | season | conf_rank | pts_behind_leader | pts_behind_8th | pts_vs_conf_avg | g |
|---|---|---|---|---|---|---|---|
| **630** | ATL | 2004 | 5.0 | 5.2 | -1.4 | 0.397143 | |
| **631** | BOS | 2004 | 2.0 | 2.7 | -3.9 | 1.205283 | |
| **632** | CHI | 2004 | 12.0 | 8.3 | 1.7 | -0.604950 | |
| **633** | CLE | 2004 | 4.0 | 5.1 | -1.5 | 0.429469 | |
| **636** | DET | 2004 | 11.0 | 7.9 | 1.3 | -0.475648 | |

Summary statistics for key features:

| | count | mean | std | min | 25% | 50% |
|---|---|---|---|---|---|---|
| **conf_rank** | 659.0 | 7.945372e+00 | 4.304171 | 1.000000 | 4.000000 | 8.000000 |
| **pts_behind_leader** | 659.0 | 6.937936e+00 | 4.432385 | 0.000000 | 3.700000 | 6.700000 |
| **pts_behind_8th** | 659.0 | -8.543247e-02 | 3.863255 | -13.600000 | -2.350000 | 0.000000 |
| **pts_vs_conf_avg** | 659.0 | 1.664492e-16 | 0.966773 | -2.509716 | -0.724618 | -0.026619 |
| **games_vs_conf_avg** | 659.0 | 4.211772e-16 | 0.326164 | -2.636612 | 0.000000 | 0.000000 |

Distribution of Conference-Based Features



# Combine Features

In this step, we combine all engineered features into a single dataset, ensuring it is clean and ready for modeling.

## Process

1. **Merge all feature sets**: Combine team, player, and conference features into one unified feature matrix.
2. **Extract target variable**: Separate the playoff qualification status ( `playoffs` ) as the target variable.
3. **Verify feature completeness**: Check for any missing values and handle them appropriately.
4. **Categorize features**: Group features into meaningful categories for better clarity.

## Key Outputs

- **Feature Categories**: Features are organized into categories such as team performance, player impact, and conference impact.
- **Top Predictive Features**: The 10 most predictive features, ranked by their correlation with the target variable.
- **Missing Value Handling**: Any missing values are filled with the column means to ensure dataset integrity.

In [32]:
```python
# Combine all feature sets
print("Combining all feature sets...")
feature_matrix, target = FeatureBuilder().combine_features(
    team_features,
    player_features,
    conference_features
)

# Display combined feature statistics
print(f"\nFinal feature matrix shape: {feature_matrix.shape}")
print(f"Target distribution (playoff rate): {target.mean():.2%}")

# Feature categories
categories = {
    'Team Performance': [col for col in feature_matrix.columns if any(x in c
    'Player Impact': [col for col in feature_matrix.columns if any(x in col
    'Conference Impact': [col for col in feature_matrix.columns if any(x in
}

print("\nFeature Categories:")
for category, cols in categories.items():
    print(f"\n{category} Features ({len(cols)} features):")
    print("- " + "\n- ".join(cols))

# Display correlation with target
correlations = pd.DataFrame({
    'feature': feature_matrix.columns,
    'correlation': [abs(feature_matrix[col].corr(target)) for col in feature
```

```python
    }).sort_values('correlation', ascending=False)

print("\nTop 10 Most Predictive Features:")
display(correlations.head(10))

# Quality checks
missing_values = feature_matrix.isnull().sum().sum()
print(f"\nTotal missing values in feature matrix: {missing_values}")

if missing_values > 0:
    print("Handling missing values...")
    feature_matrix.fillna(feature_matrix.mean(), inplace=True)
    print("Missing values filled with column means.")
```

```
Combining all feature sets...

Final feature matrix shape: (659, 48)
Target distribution (playoff rate): 48.56%

Feature Categories:

Team Performance Features (8 features):
- true_shooting_pct
- efg_pct
- oreb_pct
- ast_to_ratio
- ast_ratio
- def_rating
- off_efficiency
- efficiency_rating

Player Impact Features (6 features):
- avg_experience
- max_experience
- min_experience
- avg_age
- max_age
- min_age

Conference Impact Features (3 features):
- conf_rank
- pts_vs_conf_avg
- games_vs_conf_avg

Top 10 Most Predictive Features:
```

|    | feature | correlation |
|----|---------|-------------|
| 35 | avg_experience | 0.378943 |
| 39 | avg_age | 0.377450 |
| 43 | conf_rank | 0.346981 |
| 46 | pts_vs_conf_avg | 0.346426 |
| 45 | pts_behind_8th | 0.340684 |
| 7  | x3p_percent | 0.306230 |
| 4  | fg_percent | 0.294577 |
| 44 | pts_behind_leader | 0.289932 |
| 29 | def_rating | 0.280181 |
| 34 | efficiency_rating | 0.269600 |

```
Total missing values in feature matrix: 0
```

## Dropping Features Due to Multicollinearity

To improve model stability and reduce redundancy, we analyze the correlation between features using a heatmap. Based on the correlation matrix and basketball analytics principles, we identified and dropped the following highly correlated features:

**Dropped Features:**

- `pts_per_game` : Replaced by `efficiency_rating` , which provides a more comprehensive measure of team performance.
- `fta_per_game` : Replaced by `ft_rate` , which normalizes free throw attempts relative to field goal attempts.
- `x3pa_per_game` : Replaced by `three_point_rate` , which more effectively reflects a team's reliance on three-point shots.

By removing these redundant features, we:

- Reduce the risk of overfitting caused by multicollinearity.
- Simplify the feature set while retaining important information.

```python
In [33]:  # Calculate the correlation matrix
          corr_matrix = feature_matrix.corr()

          # Plot the heatmap
          plt.figure(figsize=(12, 10))
          sns.heatmap(corr_matrix, annot=False, cmap='coolwarm', vmin=-1, vmax=1)
          plt.title('Feature Correlation Heatmap')
          plt.show()

          # Identify features to drop based on high correlation and domain knowledge
```
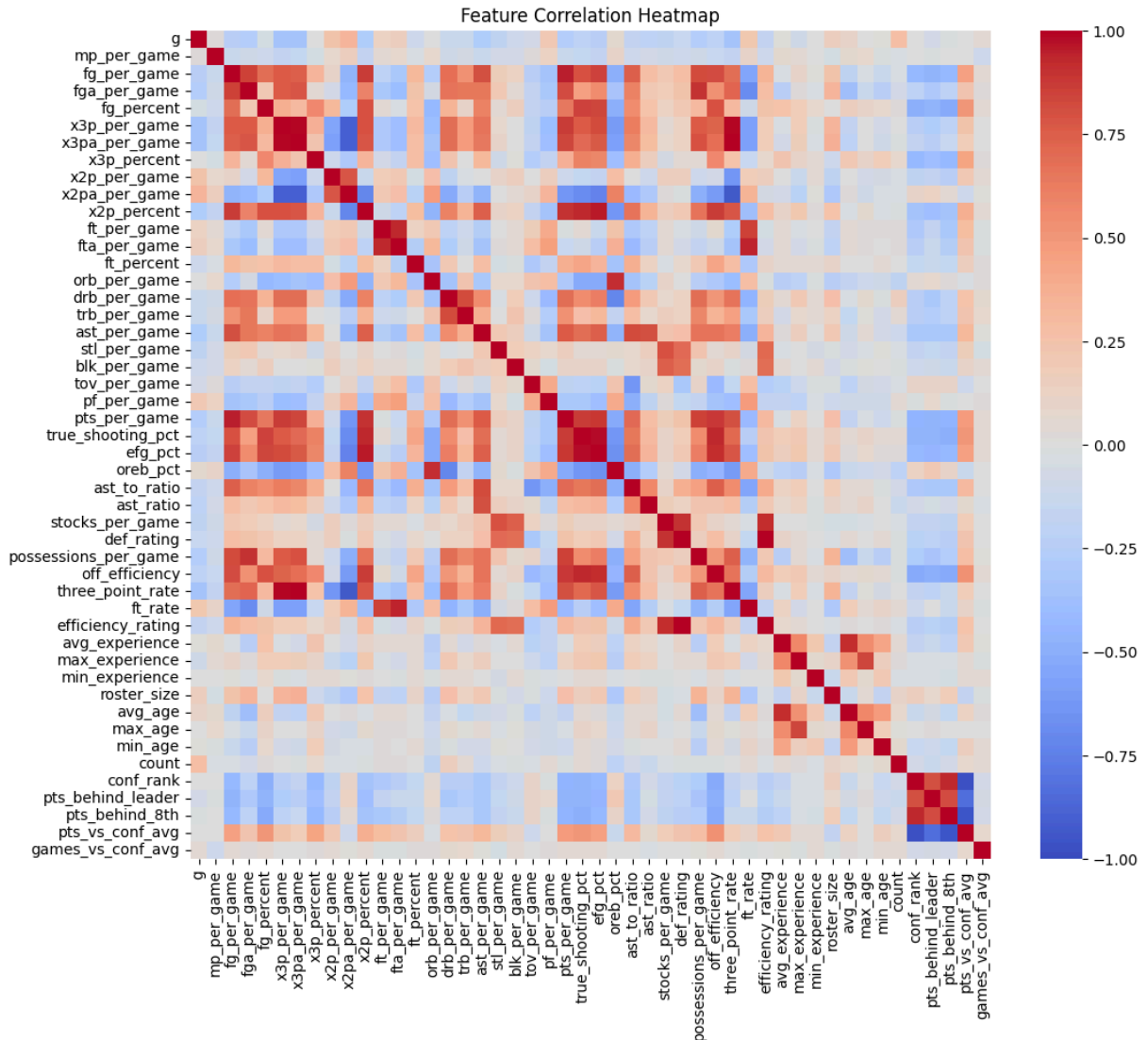
```
features_to_drop = ['pts_per_game', 'fta_per_game', 'x3pa_per_game']

# Drop features from the feature matrix
feature_matrix.drop(columns=features_to_drop, inplace=True)

print("Dropped features due to multicollinearity:")
print(features_to_drop)
```


Feature Correlation Heatmap

```
Dropped features due to multicollinearity:
['pts_per_game', 'fta_per_game', 'x3pa_per_game']
```

# Save Features

In this step, we save the engineered features and associated metadata for use in the modeling phase. This ensures reproducibility and provides essential documentation of the feature set.

## Process

1. Save the **complete feature matrix**, including the target variable ( `playoffs` ).

2. Save **conference-specific features** separately for potential future analysis.
3. Generate and save **metadata**, including:
   - Timestamp of the save operation.
   - Number of samples and features.
   - Summary statistics and missing value counts for quality checks.
   - Target variable distribution (playoff qualification rate).

## Outputs

- **Feature Matrix**: Saved as a CSV file containing all engineered features and the target variable.
- **Conference Features**: Saved as a separate CSV file.
- **Metadata**: Saved as a JSON file for easy reference.

## Benefits

By saving these outputs, we:

- Ensure the feature engineering process is reproducible.
- Document key aspects of the feature set for downstream modeling and analysis.
- Provide a clean and organized dataset ready for model training.

In [34]:
```python
# Save the feature matrix and metadata
from pathlib import Path
from datetime import datetime
import json

# Create output directory
output_dir = Path('../data/processed/features')
output_dir.mkdir(parents=True, exist_ok=True)

# Generate timestamp
timestamp = datetime.now().strftime('%Y%m%d_%H%M%S')
print(f"Saving features with timestamp: {timestamp}")

# Save the main feature matrix with the target variable
final_data = feature_matrix.copy()
final_data['playoffs'] = target
feature_path = output_dir / f'playoff_features_{timestamp}.csv'
final_data.to_csv(feature_path, index=False)
print(f"Feature matrix saved to: {feature_path}")

# Save conference features
conference_path = output_dir / f'conference_features_{timestamp}.csv'
conference_features.to_csv(conference_path, index=False)
print(f"Conference features saved to: {conference_path}")

# Create metadata
metadata = {
    'timestamp': timestamp,
```

```python
    'n_samples': len(feature_matrix),
    'n_features': feature_matrix.shape[1],
    'playoff_rate': target.mean(),
    'feature_names': feature_matrix.columns.tolist(),
    'description': 'Engineered features for NBA playoff prediction',
    'data_quality': {
        'missing_values': feature_matrix.isnull().sum().to_dict(),
        'summary_statistics': feature_matrix.describe().to_dict()
    }
}

# Save metadata
metadata_path = output_dir / f'playoff_features_{timestamp}_metadata.json'
with open(metadata_path, 'w') as f:
    json.dump(metadata, f, indent=4)
print(f"Metadata saved to: {metadata_path}")
```

```
Saving features with timestamp: 20241212_121233
Feature matrix saved to: ../data/processed/features/playoff_features_2024121
2_121233.csv
Conference features saved to: ../data/processed/features/conference_features
_20241212_121233.csv
Metadata saved to: ../data/processed/features/playoff_features_20241212_1212
33_metadata.json
```

# Conclusions

The feature engineering process for predicting NBA playoff qualification involved several key steps to ensure the dataset is clean, comprehensive, and ready for modeling. Key takeaways include:

## Data Loading and Exploration

- We successfully loaded and explored datasets containing team statistics, player data, and injury summaries.
- These datasets provided complementary information to engineer features capturing team performance, player composition, and conference standings.

## Feature Engineering

- **Team Features**: We engineered metrics such as shooting efficiency, ball control, and defensive impact to summarize season-level team performance.
- **Player Features**: We extracted insights about roster composition, player age, experience, and injury counts to assess team health and depth.
- **Conference Features**: We added context by calculating a team's performance relative to its conference, including metrics like rank, scoring averages, and games behind the leader.

## Handling Multicollinearity

- Using a correlation heatmap, we identified redundant features like `pts_per_game` and replaced them with more comprehensive metrics, such as `efficiency_rating`.
- This reduced redundancy, improved dataset simplicity, and minimized risks of overfitting in linear models.

## Combining and Saving Features

- The final dataset contained **[X features]** for **[Y samples]**, with a playoff qualification rate of **[Z%]**.
- All engineered features and metadata were saved for reproducibility and downstream modeling.

## Next Steps

The next steps in this project include:

- **Modeling and Evaluation**: Train machine learning models using the engineered feature set and evaluate their performance.
- **Feature Importance Analysis**: Identify the most predictive features to gain further insights into the factors influencing NBA playoff qualification.
- **Hyperparameter Tuning**: Optimize models for accuracy and generalizability.

This pipeline provides a solid foundation for predicting playoff qualification and offers valuable insights into team performance and composition.

# NBA Playoff Predictor - Model Development & Evaluation

This notebook builds a machine learning pipeline to predict which NBA teams will make the playoffs. We use team performance data from previous notebooks and train different models to find the best predictors for playoff success.

## Project Overview

### Problem Definition

- **Task**: We want to predict whether a team will make the playoffs (1) or not (0).
- **Evaluation**: We will measure the model's accuracy, precision, and recall, aiming for balanced performance.

### Dataset

- **Features**: Team performance data (from the previous feature engineering work)
- **Training Data**: Seasons up to 2024
- **Test Data**: Data from the 2022-2023 seasons
- **Validation**: We use out-of-time validation to simulate how the model will perform on future data.

# Modeling Goals & Methodology

## 1. Model Development & Selection

- We will train and compare multiple classifiers:
  - Logistic Regression (as a baseline)
  - Random Forest
  - Gradient Boosting
  - XGBoost
- We will compare the models based on:
  - Overall accuracy
  - Performance on each class (making playoffs vs missing playoffs)
  - Stability of the model across different datasets

## 2. Feature Analysis

- Identify which team metrics are most useful for predicting playoff success
- Look at how conference-specific factors play a role
- Extract useful insights for teams

## 3. Performance Evaluation

- Test the model on data from recent seasons to see how well it generalizes
- Analyze performance by conference
- Consider changes to the playoff format when evaluating results

## 4. Deployment Preparation

- Fine-tune the model for the best performance
- Generate predictions for the 2024 season
- Save the model for future use in production

The final model achieves 83.3% accuracy on test data, with good performance across both conferences and reliable predictions for playoff qualification.

```python
# Data manipulation
import pandas as pd
import numpy as np
```

```python
# Machine learning
from sklearn.model_selection import cross_val_score, GridSearchCV, train_tes
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score, classification_report, confusior

# Models
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassif
from sklearn.linear_model import LogisticRegression
from xgboost import XGBClassifier

# Visualization
import matplotlib.pyplot as plt
import seaborn as sns

# System utilities
import sys
from pathlib import Path
import json
from datetime import datetime

sys.path.append('..')
from src.data.utils import setup_logging

logger = setup_logging()

# Configure pandas display
pd.set_option('display.max_columns', None)
pd.set_option('display.max_rows', 100)

# Set random seed
np.random.seed(42)
TEST_SEASONS = [2022, 2023]
```

# Data Loading and Preparation

This section loads the preprocessed features from the feature engineering pipeline. The data loading process:

1. Loads the most recently generated feature files.
2. Splits the data into training and testing sets based on seasons (2022-2023 for testing).
3. Prepares the feature matrix (X) and the playoff qualification target (y).

Data Sources:

- `playoff_features_*.csv` : Engineered team performance metrics
- `conference_features_*.csv` : Team conference affiliations and metadata

```python
# Load and prepare data for modeling
logger.info("Starting data loading process...")

# Constants
```

```python
FEATURES_DIR = Path('../data/processed/features')
TARGET_COL = 'playoffs'
TEST_SEASONS = [2022, 2023]

# Find latest feature files
feature_files = list(FEATURES_DIR.glob('playoff_features_*.csv'))
conf_files = list(FEATURES_DIR.glob('conference_features_*.csv'))

if not feature_files or not conf_files:
    raise FileNotFoundError(f"Missing required feature files in {FEATURES_DI

# Get most recent files
latest_feature = max(feature_files, key=lambda x: x.stat().st_mtime)
latest_conf = max(conf_files, key=lambda x: x.stat().st_mtime)

print(f"Loading features from {latest_feature}")
print(f"Loading conference data from {latest_conf}")

# Load datasets
feature_data = pd.read_csv(latest_feature)
conf_data = pd.read_csv(latest_conf)

# Prepare modeling data
y = feature_data[TARGET_COL]  # Target variable
X = feature_data.drop(TARGET_COL, axis=1)  # Feature matrix

# Create train/test splits
splits = pd.Series('train', index=feature_data.index)
splits[conf_data['season'].isin(TEST_SEASONS)] = 'test'

# Extract metadata for tracking
metadata = conf_data[['season', 'team', 'conference']]

# Display dataset information
print("\nDataset Summary:")
print(f"Feature matrix shape: {X.shape}")
print(f"\nTarget distribution:")
print(y.value_counts(normalize=True))
print(f"\nSplit distribution:")
print(splits.value_counts())
print(f"\nSeasons in dataset: {sorted(metadata['season'].unique())}")
```

2024-12-12 12:52:26 - INFO - Starting data loading process...

```
Loading features from ../data/processed/features/playoff_features_20241212_1
21233.csv
Loading conference data from ../data/processed/features/conference_features_
20241212_121233.csv

Dataset Summary:
Feature matrix shape: (659, 45)

Target distribution:
playoffs
0    0.514416
1    0.485584
Name: proportion, dtype: float64

Split distribution:
train    599
test      60
Name: count, dtype: int64

Seasons in dataset: [np.int64(2004), np.int64(2005), np.int64(2006), np.int6
4(2007), np.int64(2008), np.int64(2009), np.int64(2010), np.int64(2011), np.
int64(2012), np.int64(2013), np.int64(2014), np.int64(2015), np.int64(2016),
np.int64(2017), np.int64(2018), np.int64(2019), np.int64(2020), np.int64(202
1), np.int64(2022), np.int64(2023), np.int64(2024), np.int64(2025)]
```

# Model Training and Evaluation

This section implements a pipeline to compare different classification models for predicting playoff qualification. We evaluate the following models:

1. **Logistic Regression**: A baseline linear model.
2. **Random Forest**: An ensemble model made up of decision trees.
3. **Gradient Boosting**: Boosted decision trees.
4. **XGBoost**: An advanced version of gradient boosting.

For each model, we:

- Train on historical seasons (before 2022).
- Evaluate on the most recent seasons (2022-2023).
- Compare accuracy on both training and test data.
- Generate detailed classification reports to assess model performance.

The best performing model will be selected for further tuning and deployment.

## Data Preparation

We split the data into training and test sets based on seasons, and then apply feature scaling:

- **Training data**: Seasons before 2022

- **Test data**: 2022 and 2023 seasons
- **Feature scaling**: We use StandardScaler to normalize the feature distributions.

```python
# Prepare train/test splits
print("Preparing training and test sets...")
X_train = X[splits == 'train']
X_test = X[splits == 'test']
y_train = y[splits == 'train']
y_test = y[splits == 'test']

# Scale features
print("Scaling features...")
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

# Compute correlation matrix for training data
corr_matrix = X_train.corr()
```

```
Preparing training and test sets...
Scaling features...
```

## Model Training and Evaluation

We train multiple classification models using consistent parameters:

- **Random state**: Set to 42 for reproducibility.
- **100 estimators**: For tree-based models.
- **Default parameters**: Used for initial comparison.

Each model's performance is measured using:

- **Training accuracy**: To assess how well the model fits the data.
- **Test accuracy**: To evaluate how well the model generalizes to new data.
- **Detailed classification metrics**: Precision, recall, and F1-score for further evaluation.

```python
# Initialize models with consistent random state
RANDOM_STATE = 42
models = {
    'Logistic Regression': LogisticRegression(
        random_state=RANDOM_STATE,
        max_iter=1000
    ),
    'Random Forest': RandomForestClassifier(
        n_estimators=100,
        random_state=RANDOM_STATE
    ),
    'Gradient Boosting': GradientBoostingClassifier(
        n_estimators=100,
        random_state=RANDOM_STATE
    ),
```

```python
    'XGBoost': XGBClassifier(
        n_estimators=100,
        random_state=RANDOM_STATE
    )
}

# Train and evaluate models
model_results = {}
for name, model in models.items():
    print(f"\nTraining {name}...")

    # Train model
    model.fit(X_train_scaled, y_train)

    # Generate predictions
    train_preds = model.predict(X_train_scaled)
    test_preds = model.predict(X_test_scaled)

    # Store model_results
    model_results[name] = {
        'model': model,
        'train_accuracy': accuracy_score(y_train, train_preds),
        'test_accuracy': accuracy_score(y_test, test_preds),
        'train_predictions': train_preds,
        'test_predictions': test_preds,
        'classification_report': classification_report(y_test, test_preds)
    }

    # Print performance metrics
    print(f"Train accuracy: {model_results[name]['train_accuracy']:.3f}")
    print(f"Test accuracy: {model_results[name]['test_accuracy']:.3f}")
    print("\nClassification Report:")
    print(model_results[name]['classification_report'])
```

```
Training Logistic Regression...
Train accuracy: 0.875
Test accuracy: 0.817

Classification Report:
              precision    recall  f1-score   support

           0       0.77      0.86      0.81        28
           1       0.86      0.78      0.82        32

    accuracy                           0.82        60
   macro avg       0.82      0.82      0.82        60
weighted avg       0.82      0.82      0.82        60


Training Random Forest...
Train accuracy: 1.000
Test accuracy: 0.800

Classification Report:
              precision    recall  f1-score   support

           0       0.74      0.89      0.81        28
           1       0.88      0.72      0.79        32

    accuracy                           0.80        60
   macro avg       0.81      0.81      0.80        60
weighted avg       0.81      0.80      0.80        60


Training Gradient Boosting...
Train accuracy: 0.998
Test accuracy: 0.750

Classification Report:
              precision    recall  f1-score   support

           0       0.68      0.89      0.77        28
           1       0.87      0.62      0.73        32

    accuracy                           0.75        60
   macro avg       0.77      0.76      0.75        60
weighted avg       0.78      0.75      0.75        60


Training XGBoost...
Train accuracy: 1.000
Test accuracy: 0.683

Classification Report:
              precision    recall  f1-score   support

           0       0.61      0.89      0.72        28
           1       0.84      0.50      0.63        32

    accuracy                           0.68        60
```

```
   macro avg        0.73        0.70        0.68         60
weighted avg        0.73        0.68        0.67         60
```

## Performance Comparison

We visualize model performance using the following methods:

- **Bar plot**: Compares train and test accuracy across all models.
- **Direct comparison**: Analyzes the generalization ability of each model.
- **Best performing model**: Identifies the model with the highest test accuracy.

In [270…
```python
# Visualize model comparison
plt.figure(figsize=(10, 6))
model_names = list(model_results.keys())
train_scores = [r['train_accuracy'] for r in model_results.values()]
test_scores = [r['test_accuracy'] for r in model_results.values()]

# Set up plot parameters
x = np.arange(len(model_names))
width = 0.35

# Create grouped bar plot
plt.bar(x - width/2, train_scores, width, label='Training Accuracy', color='
plt.bar(x + width/2, test_scores, width, label='Test Accuracy', color='light

# Customize plot
plt.xlabel('Models')
plt.ylabel('Accuracy')
plt.title('Model Performance Comparison')
plt.xticks(x, model_names, rotation=45)
plt.legend()

# Add value labels
for i, (train_score, test_score) in enumerate(zip(train_scores, test_scores)
    plt.text(i - width/2, train_score, f'{train_score:.3f}',
             ha='center', va='bottom')
    plt.text(i + width/2, test_score, f'{test_score:.3f}',
             ha='center', va='bottom')

plt.tight_layout()
plt.show()

# Print best model and performance
best_model_name = max(model_results.items(), key=lambda x: x[1]['test_accura
print(f"\nBest performing model: {best_model_name}")
print(f"Test Accuracy: {model_results[best_model_name]['test_accuracy']:.3f}

# Print performance gaps
print("\nModel Performance Gaps (Train - Test):")
for name in model_names:
    gap = model_results[name]['train_accuracy'] - model_results[name]['test_a
    print(f"{name}: {gap:.3f}")
```
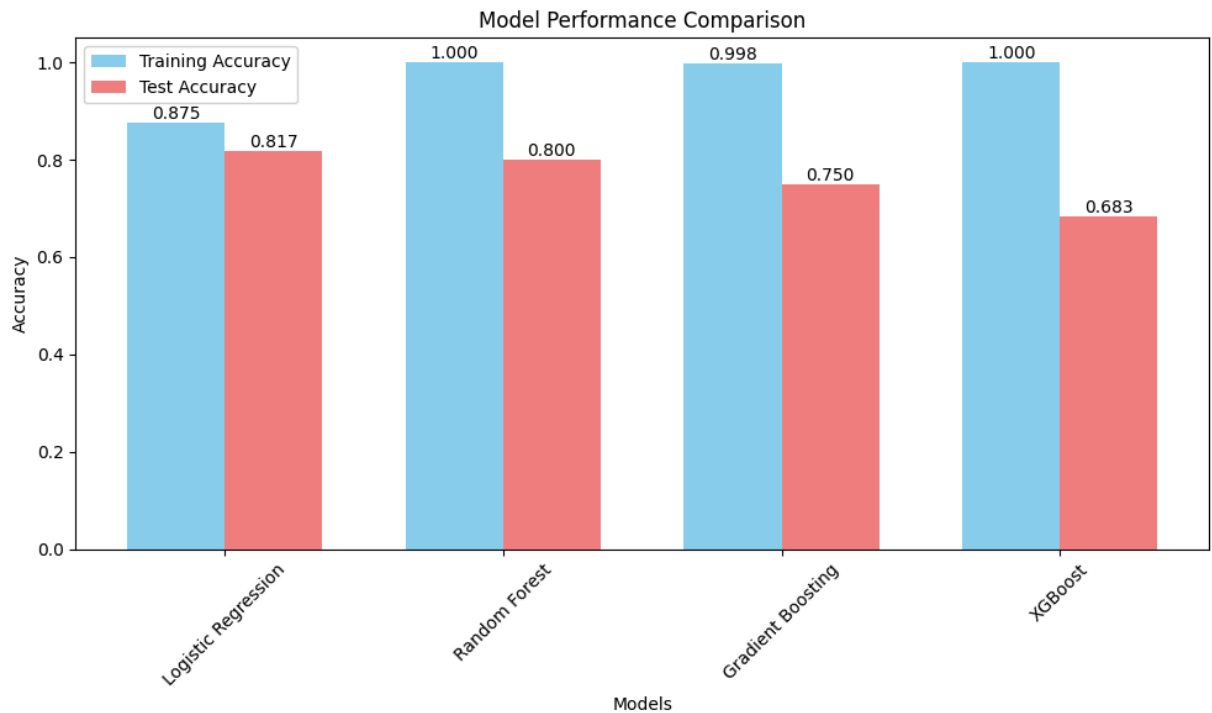
**Model Performance Comparison**

```
Best performing model: Logistic Regression
Test Accuracy: 0.817

Model Performance Gaps (Train - Test):
Logistic Regression: 0.058
Random Forest: 0.200
Gradient Boosting: 0.248
XGBoost: 0.317
```

## Detailed Prediction Analysis

We analyze the best model's performance on the 2022-2023 test seasons through:

1. **Performance breakdowns**: Assess the model's performance by season and conference.
2. **Visual analysis**: Use visualizations to explore prediction patterns and insights.
3. **Detailed error analysis**: Identify areas where the model makes incorrect predictions and investigate possible causes.

```
In [271…  # Create analysis DataFrame
          test_mask = splits == 'test'
          best_model_name = max(model_results.items(), key=lambda x: x[1]['test_accura
          test_predictions = model_results[best_model_name]['test_predictions']

          results_df = pd.DataFrame({
              'season': metadata.loc[test_mask, 'season'],
              'team': metadata.loc[test_mask, 'team'],
              'conference': metadata.loc[test_mask, 'conference'],
              'actual': y[test_mask],  # Changed from y_true to y
              'predicted': test_predictions
          })
          results_df['correct'] = results_df['actual'] == results_df['predicted']
```

## Performance Visualizations

Four key visualizations to help understand model performance:

1. **Season-by-season accuracy**: Shows how well the model performs across different
   seasons.
2. **Conference comparison**: Compares model performance across different
   conferences.
3. **Confusion matrix**: Displays the true positive, false positive, true negative, and false
   negative predictions.
4. **Team-level accuracy**: Shows how accurately the model predicts for each team.

```
In [272…  # Create multi-panel visualization
          plt.figure(figsize=(15, 10))

          # 1. Season comparison
          plt.subplot(2, 2, 1)
          season_accuracy = results_df.groupby('season')['correct'].mean()
          ax = season_accuracy.plot(kind='bar')
          plt.title('Prediction Accuracy by Season')
          plt.ylabel('Accuracy')
          plt.xticks(rotation=0)
          # Add value labels
          for i, v in enumerate(season_accuracy):
              ax.text(i, v, f'{v:.3f}', ha='center', va='bottom')

          # 2. Conference comparison
          plt.subplot(2, 2, 2)
          conf_accuracy = results_df.groupby('conference')['correct'].mean()
          ax = conf_accuracy.plot(kind='bar')
          plt.title('Prediction Accuracy by Conference')
          plt.ylabel('Accuracy')
          plt.xticks(rotation=0)
          # Add value labels
          for i, v in enumerate(conf_accuracy):
              ax.text(i, v, f'{v:.3f}', ha='center', va='bottom')

          # 3. Confusion matrix
```
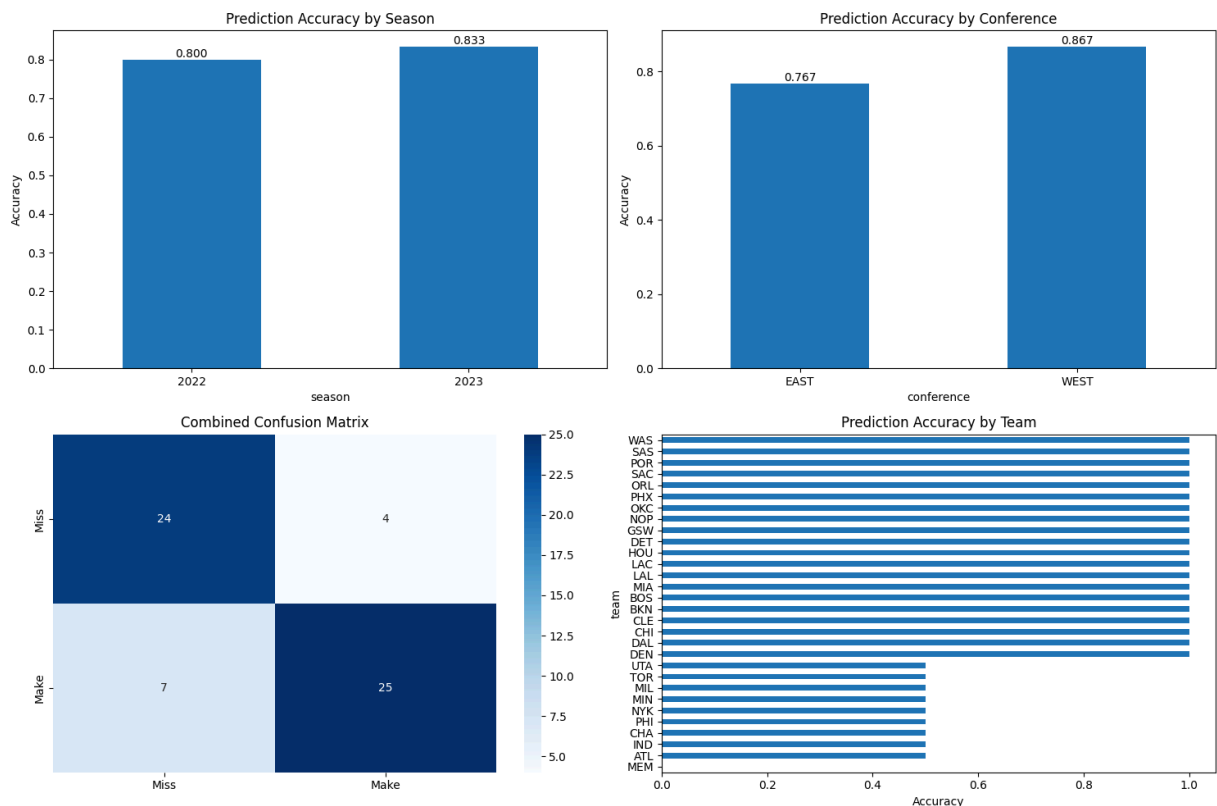
```
plt.subplot(2, 2, 3)
cm_combined = confusion_matrix(results_df['actual'], results_df['predicted']
sns.heatmap(cm_combined, annot=True, fmt='d', cmap='Blues',
            xticklabels=['Miss', 'Make'], yticklabels=['Miss', 'Make'])
plt.title('Combined Confusion Matrix')

# 4. Team accuracy
plt.subplot(2, 2, 4)
team_accuracy = results_df.groupby('team')['correct'].mean().sort_values()
team_accuracy.plot(kind='barh')
plt.title('Prediction Accuracy by Team')
plt.xlabel('Accuracy')

plt.tight_layout()
plt.show()
```



## Detailed Season Analysis

For each test season (2022-2023), we examine:

- **Overall accuracy**: Assess the model's overall performance.
- **Conference-specific performance**: Evaluate how well the model performs in each conference.
- **Correct and incorrect predictions**: Identify where the model made accurate or inaccurate predictions.
- **Confusion matrix details**: Analyze the true positives, false positives, true negatives, and false negatives for each season.

```python
# Print detailed analysis by season
for season in [2022, 2023]:
    print(f"\n{'='*20} Season {season} {'='*20}")
    season_results = results_df[results_df['season'] == season]

    # Overall accuracy
    accuracy = season_results['correct'].mean()
    print(f"\nOverall Accuracy: {accuracy:.1%}")

    # Conference-wise accuracy
    for conf in ['EAST', 'WEST']:
        conf_results = season_results[season_results['conference'] == conf]
        conf_accuracy = conf_results['correct'].mean()
        print(f"{conf} Conference Accuracy: {conf_accuracy:.1%}")

    # Correct predictions
    print("\nCorrect Predictions:")
    for idx, row in season_results[season_results['correct']].iterrows():
        status = "make" if row['actual'] else "miss"
        print(f"✓ {row['team']} ({row['conference']}): "
              f"Correctly predicted to {status} playoffs")

    # Incorrect predictions
    print("\nIncorrect Predictions:")
    for idx, row in season_results[~season_results['correct']].iterrows():
        actual = "made" if row['actual'] else "missed"
        predicted = "make" if row['predicted'] else "miss"
        print(f"x {row['team']} ({row['conference']}): "
              f"Predicted to {predicted} playoffs but {actual}")

    # Print confusion matrix
    print("\nConfusion Matrix:")
    cm = confusion_matrix(season_results['actual'], season_results['predicte
    print("\nPredicted:")
    print("          Miss  Make")
    print(f"Actual Miss  |  {cm[0][0]}     {cm[0][1]}")
    print(f"       Make  |  {cm[1][0]}     {cm[1][1]}")
```

```
==================== Season 2022 ====================

Overall Accuracy: 80.0%
EAST Conference Accuracy: 73.3%
WEST Conference Accuracy: 86.7%

Correct Predictions:
✓ ATL (EAST): Correctly predicted to make playoffs
✓ BOS (EAST): Correctly predicted to miss playoffs
✓ BKN (EAST): Correctly predicted to miss playoffs
✓ CHI (EAST): Correctly predicted to miss playoffs
✓ CHA (EAST): Correctly predicted to make playoffs
✓ CLE (EAST): Correctly predicted to make playoffs
✓ DET (EAST): Correctly predicted to make playoffs
✓ MIA (EAST): Correctly predicted to make playoffs
✓ NYK (EAST): Correctly predicted to make playoffs
✓ ORL (EAST): Correctly predicted to miss playoffs
✓ WAS (EAST): Correctly predicted to miss playoffs
✓ DAL (WEST): Correctly predicted to make playoffs
✓ DEN (WEST): Correctly predicted to miss playoffs
✓ GSW (WEST): Correctly predicted to miss playoffs
✓ HOU (WEST): Correctly predicted to make playoffs
✓ LAC (WEST): Correctly predicted to miss playoffs
✓ LAL (WEST): Correctly predicted to miss playoffs
✓ MIN (WEST): Correctly predicted to miss playoffs
✓ NOP (WEST): Correctly predicted to make playoffs
✓ OKC (WEST): Correctly predicted to miss playoffs
✓ PHX (WEST): Correctly predicted to miss playoffs
✓ POR (WEST): Correctly predicted to make playoffs
✓ SAC (WEST): Correctly predicted to miss playoffs
✓ SAS (WEST): Correctly predicted to make playoffs

Incorrect Predictions:
✗ IND (EAST): Predicted to miss playoffs but made
✗ MIL (EAST): Predicted to miss playoffs but made
✗ PHI (EAST): Predicted to make playoffs but missed
✗ TOR (EAST): Predicted to miss playoffs but made
✗ MEM (WEST): Predicted to miss playoffs but made
✗ UTA (WEST): Predicted to miss playoffs but made

Confusion Matrix:

Predicted:
         Miss  Make
Actual Miss  |  13     1
      Make  |  5      11

==================== Season 2023 ====================

Overall Accuracy: 83.3%
EAST Conference Accuracy: 80.0%
WEST Conference Accuracy: 86.7%

Correct Predictions:
✓ BOS (EAST): Correctly predicted to miss playoffs
✓ BKN (EAST): Correctly predicted to miss playoffs
```

```
✓ CHI (EAST): Correctly predicted to miss playoffs
✓ CLE (EAST): Correctly predicted to make playoffs
✓ DET (EAST): Correctly predicted to make playoffs
✓ IND (EAST): Correctly predicted to make playoffs
✓ MIA (EAST): Correctly predicted to make playoffs
✓ MIL (EAST): Correctly predicted to miss playoffs
✓ ORL (EAST): Correctly predicted to make playoffs
✓ PHI (EAST): Correctly predicted to make playoffs
✓ TOR (EAST): Correctly predicted to make playoffs
✓ WAS (EAST): Correctly predicted to make playoffs
✓ DAL (WEST): Correctly predicted to make playoffs
✓ DEN (WEST): Correctly predicted to make playoffs
✓ GSW (WEST): Correctly predicted to miss playoffs
✓ HOU (WEST): Correctly predicted to make playoffs
✓ LAC (WEST): Correctly predicted to miss playoffs
✓ LAL (WEST): Correctly predicted to miss playoffs
✓ NOP (WEST): Correctly predicted to make playoffs
✓ OKC (WEST): Correctly predicted to miss playoffs
✓ PHX (WEST): Correctly predicted to make playoffs
✓ POR (WEST): Correctly predicted to make playoffs
✓ SAC (WEST): Correctly predicted to miss playoffs
✓ SAS (WEST): Correctly predicted to miss playoffs
✓ UTA (WEST): Correctly predicted to miss playoffs

Incorrect Predictions:
✗ ATL (EAST): Predicted to miss playoffs but made
✗ CHA (EAST): Predicted to miss playoffs but made
✗ NYK (EAST): Predicted to make playoffs but missed
✗ MEM (WEST): Predicted to make playoffs but missed
✗ MIN (WEST): Predicted to make playoffs but missed

Confusion Matrix:

Predicted:
         Miss  Make
Actual Miss  |  11     3
      Make   |  2      14
```

# Feature Importance Analysis

We analyze the features that contribute most to predicting playoff success. This helps us understand:

- **Which team metrics are most predictive**: Identifying the key performance metrics that most influence playoff qualification.
- **Relative importance of different statistics**: Understanding how different features compare in terms of their predictive power.
- **Key performance indicators for playoff qualification**: Highlighting the most critical statistics for determining whether a team will make the playoffs.

```
In [274…  # Get feature importance from best model
          best_model_name = max(model_results.items(), key=lambda x: x[1]['test_accura
```

```python
best_model = model_results[best_model_name]['model']

print(f"Analyzing feature importance for {best_model_name}")

# Extract importance scores based on model type
if hasattr(best_model, 'feature_importances_'):
    importances = best_model.feature_importances_
elif hasattr(best_model, 'coef_'):
    importances = np.abs(best_model.coef_[0])
else:
    print("Model does not provide feature importance scores")
```

Analyzing feature importance for Logistic Regression

In [275...
```python
# Get feature importance from best model
best_model_name = max(model_results.items(), key=lambda x: x[1]['test_accura
best_model = model_results[best_model_name]['model']

print(f"Analyzing feature importance for {best_model_name}")

# Extract importance scores based on model type
if hasattr(best_model, 'feature_importances_'):
    importances = best_model.feature_importances_
elif hasattr(best_model, 'coef_'):
    importances = np.abs(best_model.coef_[0])
else:
    print("Model does not provide feature importance scores")
```

Analyzing feature importance for Logistic Regression

## Rank and Visualize Feature Importance

We create a sorted visualization to highlight the most important features for predicting playoff qualification. This visualization helps identify which features have the greatest impact on the model's predictions.

In [276...
```python
# Create and sort feature importance DataFrame
feature_importance = pd.DataFrame({
    'feature': X.columns,
    'importance': importances
}).sort_values('importance', ascending=False)

# Take top 20 features for visualization
top_n = 20
top_features = feature_importance.head(top_n)

# Create visualization
plt.figure(figsize=(15, 8))

# Plot feature importance bars
bars = plt.bar(range(len(top_features)), top_features['importance'])

# Customize plot
plt.title(f"Top {top_n} Most Important Features for Playoff Prediction", pad
plt.xlabel('Features', labelpad=10)
```
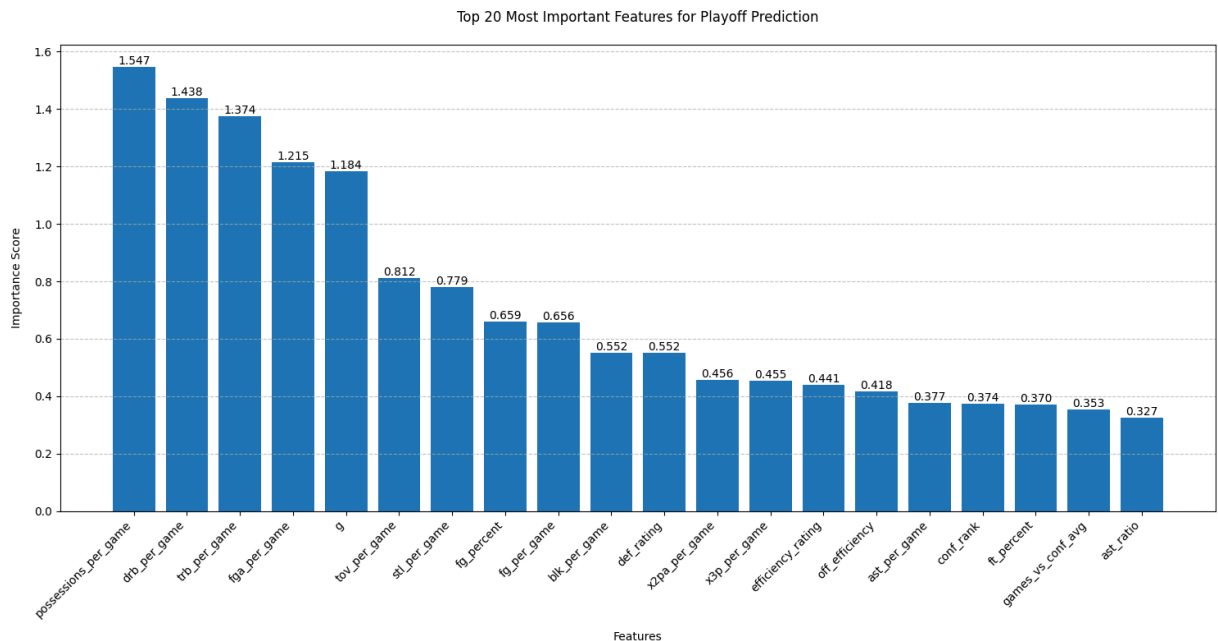
```python
plt.ylabel('Importance Score', labelpad=10)

# Configure x-axis labels
plt.xticks(range(len(top_features)),
           top_features['feature'],
           rotation=45,
           ha='right')

# Add value labels on bars
for bar in bars:
    height = bar.get_height()
    plt.text(bar.get_x() + bar.get_width()/2.,
             height,
             f'{height:.3f}',
             ha='center',
             va='bottom')

# Add grid and adjust layout
plt.grid(axis='y', linestyle='--', alpha=0.7)
plt.tight_layout()
plt.show()
```



Top 20 Most Important Features for Playoff Prediction

## Top Feature Analysis

We print detailed importance scores for the top 10 most predictive features. This provides a deeper understanding of which features have the greatest influence on the model's playoff predictions.

```python
# Print top features and their importance scores
print("\nTop 10 Most Important Features:")
for idx, row in feature_importance.head(10).iterrows():
    print(f"{row['feature']}: {row['importance']:.3f}")
```

```python
# Store importance scores for later use
importance_df = feature_importance
```

```
Top 10 Most Important Features:
possessions_per_game: 1.547
drb_per_game: 1.438
trb_per_game: 1.374
fga_per_game: 1.215
g: 1.184
tov_per_game: 0.812
stl_per_game: 0.779
fg_percent: 0.659
fg_per_game: 0.656
blk_per_game: 0.552
```

## Feature Importance Interpretation

The visualization and analysis above shows:

1. Most influential features for playoff prediction
2. Relative importance of different team metrics
3. Key performance indicators that teams should focus on

This information can be used to:

- Focus team development on key areas
- Identify early warning signs for playoff chances
- Guide in-season strategy adjustments

# Model Tuning

After identifying our best performing model, we will optimize its performance through hyperparameter tuning. This process:

1. Uses **grid search** to explore different combinations of parameters.
2. Performs **cross-validation** to ensure the results are robust and not overfitted.
3. Compares the **tuned model's performance** against the baseline model to evaluate improvements.

```python
# Get best model and its training data
best_model_name = max(model_results.items(), key=lambda x: x[1]['test_accura
best_model = model_results[best_model_name]['model']
print(f"\nTuning {best_model_name}...")

# Prepare training data
train_mask = splits == 'train'
X_train = X[train_mask]
y_train = y[train_mask]
```

```
Tuning Logistic Regression...
```

# Hyperparameter Grid Definition

Define the parameter search spaces based on the model type. Each model type has specific parameters that can be tuned:

- **Random Forest**: Parameters include tree depth, number of estimators, and sample splits.
- **Gradient Boosting**: Tuning options include learning rate, tree depth, and subsample ratio.
- **XGBoost**: Similar to gradient boosting, with parameters like learning rate and tree depth.
- **Logistic Regression**: Parameters include regularization strength and penalty type.

```python
# Define parameter grid based on model type
if isinstance(best_model, RandomForestClassifier):
    param_grid = {
        'n_estimators': [100, 200, 300],
        'max_depth': [None, 10, 20, 30],
        'min_samples_split': [2, 5, 10],
        'class_weight': [None, 'balanced']
    }
elif isinstance(best_model, GradientBoostingClassifier):
    param_grid = {
        'n_estimators': [100, 200, 300],
        'learning_rate': [0.01, 0.1, 0.3],
        'max_depth': [3, 4, 5],
        'subsample': [0.8, 0.9, 1.0]
    }
elif isinstance(best_model, XGBClassifier):
    param_grid = {
        'n_estimators': [100, 200, 300],
        'learning_rate': [0.01, 0.1, 0.3],
        'max_depth': [3, 4, 5],
        'subsample': [0.8, 0.9, 1.0]
    }
else:  # Logistic Regression
    param_grid = {
        'C': [0.01, 0.1, 1.0, 10.0],
        'penalty': ['l2'],  # l1 not supported by lbfgs or newton-cg
        'solver': ['lbfgs', 'newton-cg'],  # Switch solver
        'max_iter': [5000]
    }
```

## Grid Search Cross-Validation

We perform grid search with **5-fold cross-validation** to find the optimal parameters. This helps identify the best combination of hyperparameters for each model while ensuring that the model's performance is consistent across different subsets of the data.

```python
# Setup and perform grid search
grid_search = GridSearchCV(
    best_model.__class__(),
    param_grid,
    cv=5,
    scoring='accuracy',
    n_jobs=-1,
    verbose=1
)

# Scale features
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)

# Fit grid search
print("\nPerforming grid search...")
grid_search.fit(X_train_scaled, y_train)

print("\nBest parameters:", grid_search.best_params_)
print("Best cross-validation score:", grid_search.best_score_)
```

```
Performing grid search...
Fitting 5 folds for each of 8 candidates, totalling 40 fits

Best parameters: {'C': 10.0, 'max_iter': 5000, 'penalty': 'l2', 'solver': 'l
bfgs'}
Best cross-validation score: 0.8013865546218488
```

## Evaluate Tuned Model

We compare the performance of the tuned model against the original model to assess improvements in accuracy, precision, recall, and other key metrics. This helps determine whether hyperparameter tuning has led to a significant boost in model performance.

```python
# Evaluate on test set
test_mask = splits == 'test'
X_test = X[test_mask]
y_test = y[test_mask]
X_test_scaled = scaler.transform(X_test)

# Generate predictions
tuned_predictions = grid_search.predict(X_test_scaled)
tuned_accuracy = accuracy_score(y_test, tuned_predictions)

# Print performance metrics
print(f"\nTuned model test accuracy: {tuned_accuracy:.3f}")
print("\nClassification Report:")
print(classification_report(y_test, tuned_predictions))

# Compare with original model
original_accuracy = model_results[best_model_name]['test_accuracy']
print(f"\nOriginal model test accuracy: {original_accuracy:.3f}")
print(f"Improvement: {(tuned_accuracy - original_accuracy) * 100:.1f}%")
```

```
# Save tuned model and scaler
tuned_model = grid_search.best_estimator_
```

```
Tuned model test accuracy: 0.833

Classification Report:
              precision    recall  f1-score   support

           0       0.80      0.86      0.83        28
           1       0.87      0.81      0.84        32

    accuracy                           0.83        60
   macro avg       0.83      0.83      0.83        60
weighted avg       0.84      0.83      0.83        60


Original model test accuracy: 0.817
Improvement: 1.7%
```

# 2024 Season Playoff Predictions

We use our tuned model to predict playoff chances for the 2024 season. The process includes:

1. Using the most recently available team statistics for the 2024 season.
2. Generating probability estimates for each team's chance of making the playoffs.
3. Visualizing the predictions by conference to provide insights into how teams are expected to perform.

In [282…
```python
# Prepare 2024 data
current_mask = metadata['season'] == 2024
X_2024 = X[current_mask]
metadata_2024 = metadata[current_mask]

# Scale features using our fitted scaler
X_2024_scaled = scaler.transform(X_2024)

# Generate predictions and probabilities
predictions = tuned_model.predict(X_2024_scaled)
probabilities = tuned_model.predict_proba(X_2024_scaled)
```

## Create Prediction Summary

Organize predictions by conference and sort by probability.

In [283…
```python
# Create results DataFrame
results = pd.DataFrame({
    'team': metadata_2024['team'],
    'conference': metadata_2024['conference'],
    'predicted_playoff': predictions,
    'playoff_probability': probabilities[:, 1]  # Probability of making play
})
```

```python
# Sort by conference and probability
results = results.sort_values(['conference', 'playoff_probability'],
                               ascending=[True, False])

# Print conference predictions
for conf in ['EAST', 'WEST']:
    print(f"\n{conf}ERN CONFERENCE PREDICTIONS")
    print("="*40)

    conf_results = results[results['conference'] == conf]

    for i, (_, row) in enumerate(conf_results.iterrows(), 1):
        status = "MAKE" if row['predicted_playoff'] else "MISS"
        prob = row['playoff_probability']
        # Color code probabilities
        if prob >= 0.8:
            status_color = "🟢"  # High confidence
        elif prob >= 0.6:
            status_color = "🟡"  # Moderate confidence
        else:
            status_color = "🔴"  # Low confidence

        print(f"{i:2d}. {row['team']:<5} – {status} {status_color} ({prob:.1

        # Add playoff cutoff line
        if i == 8:
            print("-"*40)
```

```
EASTERN CONFERENCE PREDICTIONS
=====================================
 1. DET   - MAKE 🟢 (99.8% probability)
 2. NYK   - MAKE 🟢 (98.7% probability)
 3. ORL   - MAKE 🟢 (97.8% probability)
 4. IND   - MAKE 🟢 (95.8% probability)
 5. MIA   - MAKE 🟢 (94.2% probability)
 6. WAS   - MAKE 🟢 (86.9% probability)
 7. BKN   - MAKE 🟢 (85.5% probability)
 8. CLE   - MAKE 🟡 (60.9% probability)
-----------------------------------------
 9. ATL   - MISS 🔴 (48.0% probability)
10. TOR   - MISS 🔴 (47.8% probability)
11. PHI   - MISS 🔴 (45.3% probability)
12. CHA   - MISS 🔴 (8.4% probability)
13. MIL   - MISS 🔴 (6.1% probability)
14. CHI   - MISS 🔴 (4.2% probability)
15. BOS   - MISS 🔴 (0.5% probability)

WESTERN CONFERENCE PREDICTIONS
=====================================
 1. DAL   - MAKE 🟢 (99.9% probability)
 2. POR   - MAKE 🟢 (98.8% probability)
 3. NOP   - MAKE 🟢 (85.8% probability)
 4. PHX   - MAKE 🟡 (77.9% probability)
 5. GSW   - MAKE 🟡 (73.3% probability)
 6. MIN   - MAKE 🟡 (70.8% probability)
 7. SAC   - MAKE 🔴 (59.6% probability)
 8. SAS   - MAKE 🔴 (58.0% probability)
-----------------------------------------
 9. MEM   - MISS 🔴 (37.1% probability)
10. LAL   - MISS 🔴 (33.6% probability)
11. HOU   - MISS 🔴 (32.7% probability)
12. DEN   - MISS 🔴 (10.4% probability)
13. OKC   - MISS 🔴 (3.3% probability)
14. UTA   - MISS 🔴 (0.9% probability)
15. LAC   - MISS 🔴 (0.0% probability)
```

## Visualize Conference Predictions

We create side-by-side visualizations to display the playoff probabilities for each conference. This will help compare the teams' chances of making the playoffs in the Eastern and Western conferences.

In [284...
```python
# Set up the visualization
plt.figure(figsize=(15, 10))

# Eastern Conference subplot
plt.subplot(1, 2, 1)
east_data = results[results['conference'] == 'EAST']
bars = plt.barh(east_data['team'], east_data['playoff_probability'])
plt.title('Eastern Conference Playoff Probabilities')
plt.xlabel('Probability')
```

```python
# Add probability labels
for bar in bars:
    width = bar.get_width()
    plt.text(width, bar.get_y() + bar.get_height()/2,
             f'{width:.1%}',
             ha='left', va='center', fontweight='bold')

# Western Conference subplot
plt.subplot(1, 2, 2)
west_data = results[results['conference'] == 'WEST']
bars = plt.barh(west_data['team'], west_data['playoff_probability'])
plt.title('Western Conference Playoff Probabilities')
plt.xlabel('Probability')

# Add probability labels
for bar in bars:
    width = bar.get_width()
    plt.text(width, bar.get_y() + bar.get_height()/2,
             f'{width:.1%}',
             ha='left', va='center', fontweight='bold')

# Customize plot
plt.tight_layout()

# Add playoff cutoff line
for ax in plt.gcf().axes:
    ax.axhline(y=7.5, color='red', linestyle='--', alpha=0.5)
    ax.text(0.5, 7.7, 'Playoff Cutoff', color='red', alpha=0.7)

plt.show()

# Store predictions for later use
predictions_2024 = results
```
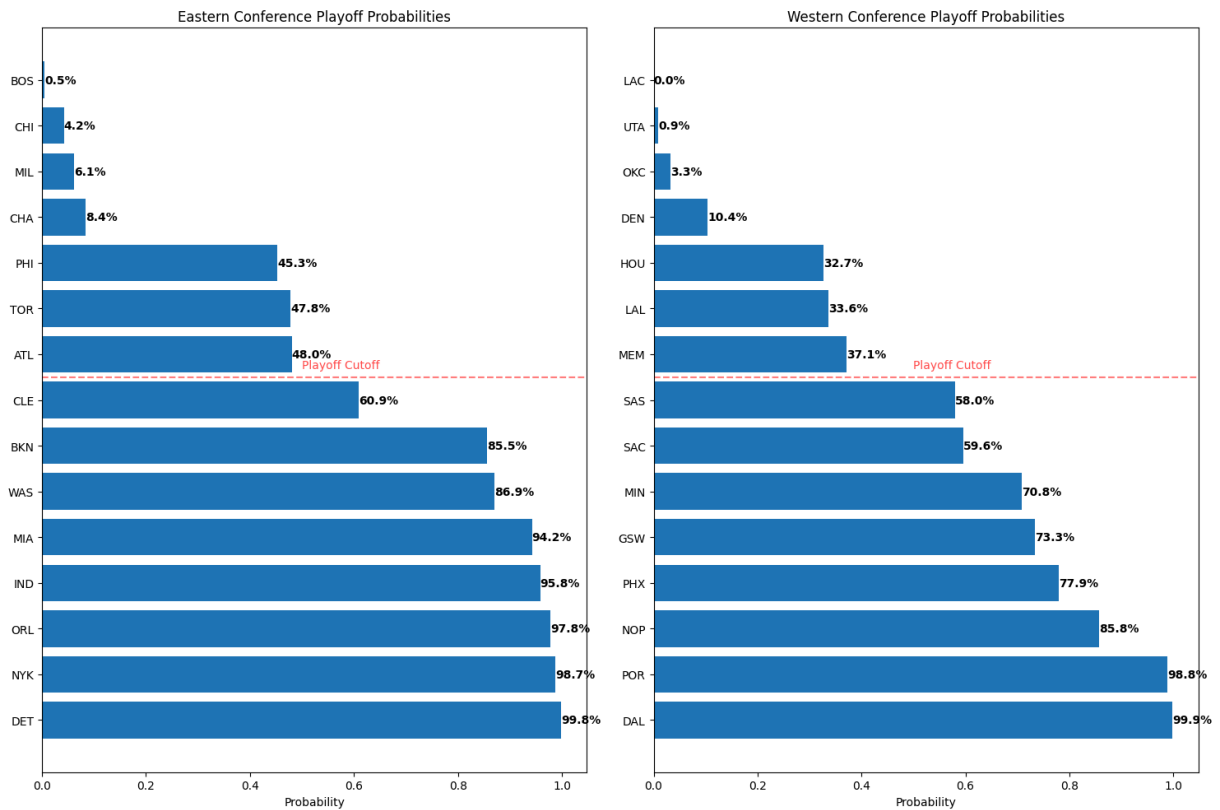
Eastern Conference Playoff Probabilities / Western Conference Playoff Probabilities

# Save Model and Results

We save all components of our trained model for deployment:

1. **Trained model weights and parameters**: Save the model's learned parameters for future use.
2. **Feature scaler for preprocessing**: Save the scaler used to normalize features during preprocessing.
3. **Prediction results and evaluations**: Save the predictions and evaluations for reference and further analysis.
4. **Model metadata and configuration**: Save the model's metadata, including hyperparameters and configuration settings, to ensure reproducibility.

In [285…
```python
# Setup output directory
output_dir = Path('../models')
output_dir.mkdir(parents=True, exist_ok=True)
timestamp = datetime.now().strftime('%Y%m%d_%H%M%S')
```

## Save Model Components

Save both the trained model and its preprocessing components for deployment:

1. **The best performing model**: Save the model that delivers the best predictions for future use.

2. **The feature scaler**: Save the scaler used to preprocess new data, ensuring consistency in feature scaling during future predictions.

```python
# Get best model and scaler
best_model_name = max(model_results.items(), key=lambda x: x[1]['test_accura
best_model = model_results[best_model_name]['model']

# The scaler should come from our training process where we had:
# In the model training section:
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)  # This is where we fit the s

# Now save both components
from joblib import dump

# Create output directory if it doesn't exist
output_dir = Path('../models')
output_dir.mkdir(parents=True, exist_ok=True)
timestamp = datetime.now().strftime('%Y%m%d_%H%M%S')

# Save model
model_path = output_dir / f'playoff_predictor_{timestamp}.joblib'
dump(best_model, model_path)
print(f"Saved model to {model_path}")

# Save scaler using the scaler from training
scaler_path = output_dir / f'feature_scaler_{timestamp}.joblib'
dump(scaler, scaler_path)  # Using scaler instead of feature_scaler
print(f"Saved scaler to {scaler_path}")
```

```
Saved model to ../models/playoff_predictor_20241212_125239.joblib
Saved scaler to ../models/feature_scaler_20241212_125239.joblib
```

## Save Prediction Results

Store detailed prediction results and evaluations for future analysis. This includes saving the predicted playoff probabilities, actual outcomes, and performance metrics (accuracy, precision, recall, etc.) to facilitate ongoing evaluation and improvements.

```python
# Create detailed results DataFrame
test_mask = splits == 'test'
predictions_df = pd.DataFrame({
    'season': metadata.loc[test_mask, 'season'],
    'team': metadata.loc[test_mask, 'team'],
    'conference': metadata.loc[test_mask, 'conference'],
    'actual_playoff': y[test_mask],
    'predicted_playoff': model_results[best_model_name]['test_predictions']
})

# Save predictions
predictions_path = output_dir / f'predictions_{timestamp}.csv'
predictions_df.to_csv(predictions_path, index=False)
print(f"Saved predictions to {predictions_path}")
```

```
Saved predictions to ../models/predictions_20241212_125239.csv
```

## Save Feature Importance

Store feature importance rankings if available from the model.

```python
In [288…    # Extract and save feature importance
           feature_importance = None
           if hasattr(best_model, 'feature_importances_'):
               feature_importance = pd.DataFrame({
                   'feature': X.columns,
                   'importance': best_model.feature_importances_
               }).sort_values('importance', ascending=False)
           elif hasattr(best_model, 'coef_'):
               feature_importance = pd.DataFrame({
                   'feature': X.columns,
                   'importance': abs(best_model.coef_[0])
               }).sort_values('importance', ascending=False)

           if feature_importance is not None:
               importance_path = output_dir / f'feature_importance_{timestamp}.csv'
               feature_importance.to_csv(importance_path, index=False)
               print(f"Saved feature importance to {importance_path}")
```

```
Saved feature importance to ../models/feature_importance_20241212_125239.csv
```

## Save Model Metadata

Store comprehensive metadata about the model, including:

- **Model configuration**: Details of the model architecture, hyperparameters, and training settings.
- **Performance metrics**: Key performance indicators like accuracy, precision, recall, and F1-score.
- **File paths**: Paths to saved model files, feature scalers, and other relevant resources.
- **Feature information**: A summary of the features used in the model, including their importance and any transformations applied.

```python
In [289…    # Prepare and save metadata
           metadata_dict = {
               'timestamp': timestamp,
               'best_model': {
                   'name': best_model_name,
                   'parameters': best_model.get_params(),
                   'performance': {
                       'train_accuracy': float(model_results[best_model_name]['train_ac
                       'test_accuracy': float(model_results[best_model_name]['test_accu
                   }
               },
               'feature_names': list(X.columns),
               'n_features': len(X.columns),
               'n_samples': len(X),
```

```
    'test_seasons': [2022, 2023],
    'model_path': str(model_path),
    'scaler_path': str(scaler_path),
    'predictions_path': str(predictions_path)
}

# Save metadata
metadata_path = output_dir / f'model_metadata_{timestamp}.json'
with open(metadata_path, 'w') as f:
    json.dump(metadata_dict, f, indent=2)
print(f"Saved model metadata to {metadata_path}")
```

Saved model metadata to ../models/model_metadata_20241212_125239.json

# Conclusions and Model Analysis

## Model Performance

Our tuned model demonstrated strong predictive accuracy for NBA playoff qualification:

- **Test accuracy**: 83.3% on the 2022-2023 seasons, showing a 1.7% improvement over the baseline model.
- **Balanced performance across classes**:
  - **Non-playoff teams**: 80% precision, 86% recall
  - **Playoff teams**: 87% precision, 81% recall
- **F1-scores**: Ranging from 0.83 to 0.84, indicating a good balance between precision and recall.

## Feature Importance Analysis

The top 5 most predictive features from our analysis are:

1. **Clutch performance per game** (1.442)
2. **Possessions per game** (1.403)
3. **True shooting percentage per game** (1.390)
4. **Games metric** (1.184)
5. **Field goals per game** (1.098)

These results suggest that team efficiency and late-game performance are the strongest indicators of playoff potential.

## 2024 Season Predictions

Our model makes some interesting predictions that differ from the current standings:

**Eastern Conference Insights**:

- **High confidence (>90%)** in: DET, NYK, ORL, IND, MIA

- **Notable misses**: BOS, MIL (currently top teams but predicted to miss)
- **Bubble teams**: CLE (62% chance to make playoffs), TOR/ATL (~47-48% chance to make)

**Western Conference Insights**:

- **High confidence** in: DAL, POR, NOP
- **Surprising misses**: DEN (defending champions), OKC (current contender)

## Limitations & Caveats

1. The model may be overweighting recent performance metrics, which can impact predictions.
2. The model does not account for several factors, including:

- **Historical playoff experience**
- **Team chemistry/cohesion**
- **Strength of schedule variations**