



TECHNISCHE
UNIVERSITÄT
WIEN

Data management plan (DMP)

Home Advantage in Professional Football – Research Data Management Project

Home Advantage in Football

Version	Effective date	Description of document/changes
1.0	22/11/2025	First version of the DMP – created for the start of the project
1.1	23/11/2025	Updated data sources section to include detailed information about reused datasets (R1, R2)
1.2	25/11/2025	Added produced datasets table and refined documentation

Level of
distribution



This DMP is licensed under a [Creative Commons Attribution 4.0 International License \(CC BY 4.0\)](#).

This DMP is not published in a public repository and does not have a DOI.

FWF Data Management Plan (DMP)

I General Information																																		
I.1 Administrative information	PI: Bartosz Ciesielski, bartosz.ciesielski@student.tuwien.ac.at, TU Wien, ROR: ror.org/04d836q62, Project Leader FWF project number: Internal Project ID: DMP version: 01, 22.11.2025 Contributors: Tomasz Miksa, tomasz.miksa@tuwien.ac.at, ORCID: 0000-0002-4929-7875, TU Wien, ROR: ror.org/04d836q62, Hosting Institution																																	
	Person responsible for data management and DMP: Bartosz Ciesielski, bartosz.ciesielski@student.tuwien.ac.at, TU Wien, ROR: ror.org/04d836q62 Co-ordination of data management responsibilities across partners: Bartosz Ciesielski, bartosz.ciesielski@student.tuwien.ac.at, TU Wien, ROR: ror.org/04d836q62 Resources costed in for RDM: There are no costs dedicated to data management and ensuring that data will be FAIR.																																	
	II Data Characteristics																																	
II.1 Data description and collection or re-use of existing data	Produced datasets: <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="text-align: left; padding: 5px;">dataset ID</th><th style="text-align: left; padding: 5px;">title</th><th style="text-align: left; padding: 5px;">type</th><th style="text-align: left; padding: 5px;">format</th><th style="text-align: left; padding: 5px;">estimated volume</th><th style="text-align: left; padding: 5px;">contains sensitive data</th><th style="text-align: left; padding: 5px;">description</th></tr> </thead> <tbody> <tr> <td style="padding: 5px;">P1</td><td style="padding: 5px;">all_matches</td><td style="padding: 5px;">Tabular data</td><td style="padding: 5px;">CSV/XLSX</td><td style="padding: 5px;">~50–100 KB</td><td style="padding: 5px;">no</td><td style="padding: 5px;">Full combined match dataset with basic preprocessing applied</td></tr> <tr> <td style="padding: 5px;">P2</td><td style="padding: 5px;">league_results_counts</td><td style="padding: 5px;">Tabular data</td><td style="padding: 5px;">CSV/XLSX</td><td style="padding: 5px;">~50–100 KB</td><td style="padding: 5px;">no</td><td style="padding: 5px;">Match result counts grouped by league</td></tr> <tr> <td style="padding: 5px;">P3</td><td style="padding: 5px;">league_results_percent</td><td style="padding: 5px;">Tabular data</td><td style="padding: 5px;">CSV/XLSX</td><td style="padding: 5px;">~50–100 KB</td><td style="padding: 5px;">no</td><td style="padding: 5px;">Match result percentages grouped by league</td></tr> </tbody> </table>						dataset ID	title	type	format	estimated volume	contains sensitive data	description	P1	all_matches	Tabular data	CSV/XLSX	~50–100 KB	no	Full combined match dataset with basic preprocessing applied	P2	league_results_counts	Tabular data	CSV/XLSX	~50–100 KB	no	Match result counts grouped by league	P3	league_results_percent	Tabular data	CSV/XLSX	~50–100 KB	no	Match result percentages grouped by league
dataset ID	title	type	format	estimated volume	contains sensitive data	description																												
P1	all_matches	Tabular data	CSV/XLSX	~50–100 KB	no	Full combined match dataset with basic preprocessing applied																												
P2	league_results_counts	Tabular data	CSV/XLSX	~50–100 KB	no	Match result counts grouped by league																												
P3	league_results_percent	Tabular data	CSV/XLSX	~50–100 KB	no	Match result percentages grouped by league																												

	P4	overall_distribution	Tabular data	CSV/XLSX	~50–100 KB	no	Overall distribution of match outcomes across both leagues																								
	P5	summary_all_matches	Tabular data	CSV/XLSX	~50–100 KB	no	Summary statistics for all matches combined																								
	P6	summary_by_league	Tabular data	CSV/XLSX	~50–100 KB	no	Summary statistics grouped by league																								
Reused datasets:																															
<table border="1"> <thead> <tr> <th>dataset ID</th><th>title</th><th>source</th><th>rights (e.g. license)</th><th>contains sensitive data</th><th>description</th><th></th><th></th></tr> </thead> <tbody> <tr> <td>R1</td><td>Professional football match results (Ekstraklasa)</td><td>https://www.football-data.co.uk/poland.php</td><td></td><td>no</td><td>...</td><td></td><td></td></tr> <tr> <td>R2</td><td>Professional football match results (Premier League)</td><td>https://github.com/jalapic/engsoccerdata/tree/master/data</td><td></td><td>no</td><td>...</td><td></td><td></td></tr> </tbody> </table>								dataset ID	title	source	rights (e.g. license)	contains sensitive data	description			R1	Professional football match results (Ekstraklasa)	https://www.football-data.co.uk/poland.php		no	...			R2	Professional football match results (Premier League)	https://github.com/jalapic/engsoccerdata/tree/master/data		no	...		
dataset ID	title	source	rights (e.g. license)	contains sensitive data	description																										
R1	Professional football match results (Ekstraklasa)	https://www.football-data.co.uk/poland.php		no	...																										
R2	Professional football match results (Premier League)	https://github.com/jalapic/engsoccerdata/tree/master/data		no	...																										
Methods and software used for data generation and reuse																															
No new data will be generated. Existing football match datasets will be reused and processed using Python (pandas) and Jupyter Notebooks. Raw data will be stored in /data/raw/ and processed outputs in /data/processed/.																															
III Documentation and Data Quality																															
III.1 Metadata and documentation	Data organisation, metadata, and documentation:																														
	The project follows a clear folder structure separating raw data (/data/raw/), processed data (/data/processed/), notebooks (/notebooks/), reports (/reports/), and source code (/src/). Raw data will remain unchanged, while all transformations and analyses will be stored as new processed files to ensure traceability. Versioning is handled through Git and GitHub. All changes to data processing scripts, notebooks, and documentation are tracked automatically via commits, providing a complete and transparent version history. File naming will follow a consistent structure that reflects content and processing steps. This setup ensures reproducibility, clear documentation of changes, and long-term maintainability of the research data.																														
The project will provide descriptive metadata to ensure that the datasets can be easily identified, understood, and reused. Metadata will include: dataset title, source URL, data origin, time period covered, variable descriptions (teams, goals, match results, home/away indicators), file formats, and licensing																															

	<p>information. As there are no specific domain metadata standards required for this type of sports statistics data, metadata will be documented at the project level through a detailed README file. Each dataset will include clear naming conventions and accompanying descriptions in the /data/processed/ directory. Information about data provenance, reuse conditions, and processing steps will be provided in the documentation so that users can trace how the raw data were transformed. This ensures transparency, discoverability, and reusability of all research outputs. This will help others to identify, discover and reuse our data.</p> <p>Additionally, we will provide common metadata such as title, description, or keywords when publishing data in open access repositories. In such a case, we will follow the default template provided by the repository, such as Data Cite Metadata or Dublin Core.</p> <p>As far as possible, we will use controlled vocabularies for our data to allow inter-disciplinary interoperability and machine-actionability.</p> <p>All analysis steps are documented in the Jupyter notebook (home_advantage_study.ipynb), including data loading, cleaning, transformations, and visualizations. The repository README describes the project structure, datasets, processing workflow, software used, and instructions for reproducing results. Processed data files and visual outputs are stored in /data/processed/ and /reports/figures/ to enable verification and reuse.</p>												
III.2 Data quality control	<p>Data quality control:</p> <p>The following data quality checks will be done: peer review of data.</p>												
IV Data Storage, Sharing, and Long-Term Preservation													
IV.1 Data storage and backup during the research process	<p>Storage and backup facilities:</p> <p>No storage options are specified at the moment.</p> <p>Data security and protection of sensitive data:</p> <p>We pay strict attention to compliance with the relevant institutional and national data protection policies. At this stage, it is not foreseen to process any sensitive data in the project. If this changes, advice will be sought from the data protection specialist at TU Wien, and the DMP will be updated.</p> <p>Access to data during research:</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="text-align: center; padding: 5px;">dataset ID</th> <th style="text-align: center; padding: 5px;">selected project members</th> <th style="text-align: center; padding: 5px;">all other project members</th> <th style="text-align: center; padding: 5px;">the public</th> </tr> </thead> <tbody> <tr> <td style="text-align: center; padding: 5px;">R1</td> <td style="text-align: center; padding: 5px;">writing</td> <td style="text-align: center; padding: 5px;">reading only</td> <td style="text-align: center; padding: 5px;">no access</td> </tr> <tr> <td style="text-align: center; padding: 5px;">R2</td> <td style="text-align: center; padding: 5px;">writing</td> <td style="text-align: center; padding: 5px;">reading only</td> <td style="text-align: center; padding: 5px;">no access</td> </tr> </tbody> </table>	dataset ID	selected project members	all other project members	the public	R1	writing	reading only	no access	R2	writing	reading only	no access
dataset ID	selected project members	all other project members	the public										
R1	writing	reading only	no access										
R2	writing	reading only	no access										

	<p>Data publication and access conditions:</p> <p>As far as possible, obtained datasets will be published in repositories. Details on access conditions, reuse licenses, reasons for restrictions, etc. are collected in the table below.</p> <table border="1"> <thead> <tr> <th>dataset ID</th><th>access conditions</th><th>estimated publication date</th><th>location for publication (repository)</th><th>PID</th><th>license</th></tr> </thead> <tbody> <tr> <td>-</td><td>-</td><td>-</td><td>-</td><td>-</td><td>-</td></tr> </tbody> </table> <p>No new datasets will be published. All reused datasets remain under the licenses of their original providers.</p>	dataset ID	access conditions	estimated publication date	location for publication (repository)	PID	license	-	-	-	-	-	-
dataset ID	access conditions	estimated publication date	location for publication (repository)	PID	license								
-	-	-	-	-	-								
IV.2 Data sharing and long-term preservation	<p>Methods or software needed to access and use data: The data are provided in standard tabular formats (.xlsx and .csv), which can be opened with common tools such as Excel, Python (pandas), or R. No specialized software is required. The folder structure and documentation in the repository support straightforward reuse.</p> <p>Long-term preservation and deletion of data:</p> <table border="1"> <thead> <tr> <th>dataset ID</th><th>location for long-term storage</th><th>minimum retention period (≥ 10 years)</th><th>foreseeable research uses and/or users</th></tr> </thead> <tbody> <tr> <td>-</td><td>-</td><td>-</td><td>-</td></tr> </tbody> </table> <p>No new datasets will be published. The project uses only existing openly accessible datasets (R1 and R2), which remain available under their original licenses. No additional repository deposit is required.</p>	dataset ID	location for long-term storage	minimum retention period (≥ 10 years)	foreseeable research uses and/or users	-	-	-	-				
dataset ID	location for long-term storage	minimum retention period (≥ 10 years)	foreseeable research uses and/or users										
-	-	-	-										
V Legal and Ethical Aspects													
V.1 Legal aspects	<p>Personal data:</p> <p>At this stage, it is not foreseen to process any personal data in the project. If this changes, advice will be sought from the data protection specialist at TU Wien, and the DMP will be updated.</p> <p>Intellectual property rights and rights of use:</p> <p>The following individual(s) hold rights and control access to the project data: Access to all datasets will be controlled solely by the project owner (Bartosz Ciesielski). No additional project members have access rights.</p>												

V.2 Ethical aspects	<p>Ethical issues: No particular ethical issue is foreseen with the data to be used or produced by the project. This section will be updated if issues arise.</p>
----------------------------	--