# Data Warehousing Laboratory Project following the CRISP-DM through the example
# of soccer statistics

**Students:** Onur Yavuz, Eugenio Donaque, Artur Baliet, Hannes Daniel

# **Structure**

1. Introduction

2. CRISP-DM

3. Concept of the problem

4. Data Preparation

    1. Solution approach

    2. Implementation steps

5. Conclusion

6. References

# **Introduction**

- Football is the most popular sport on this planet and games from previous seasons are compared to follow the development

- Our task was to enable useful information out of the huge amount of data

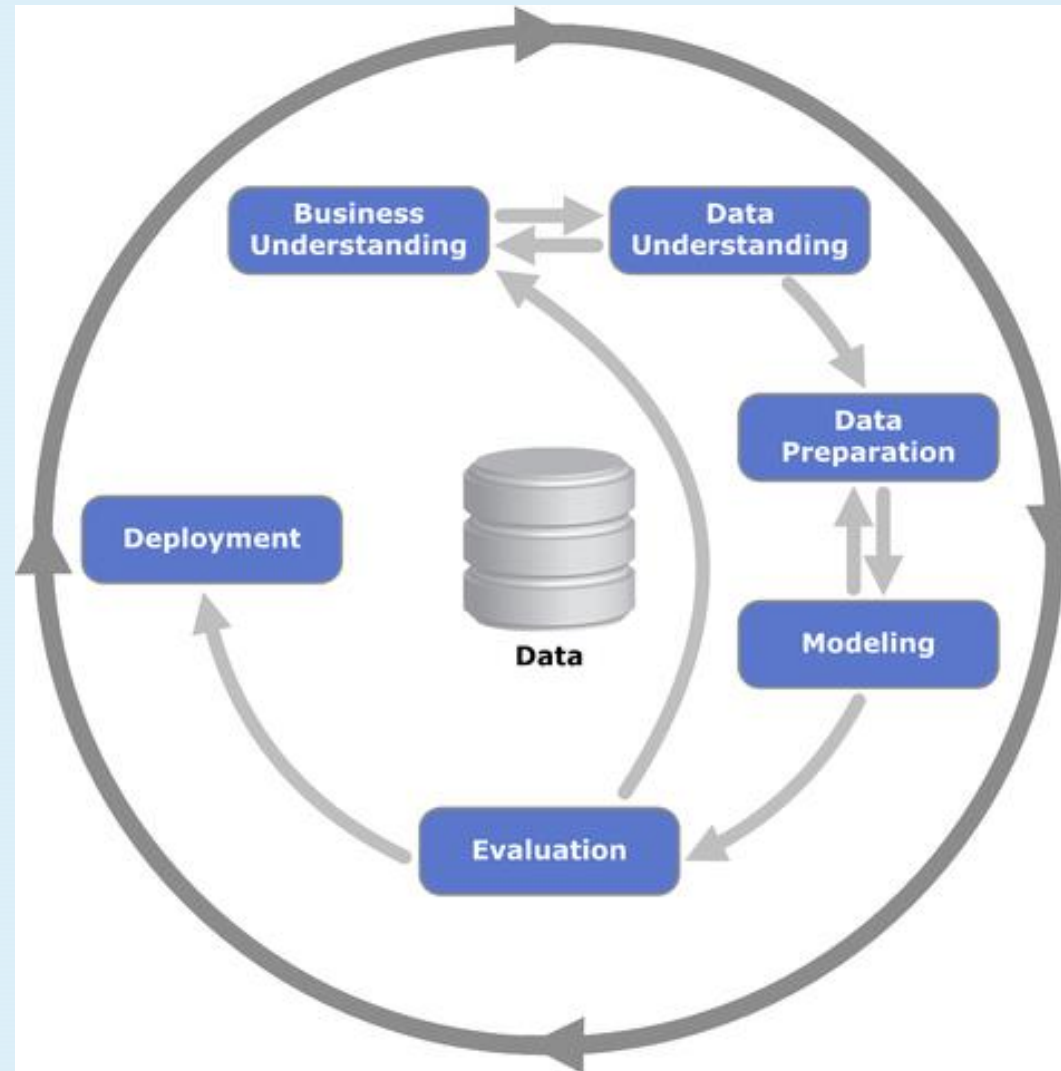- We solved this task through following the CRISP-DM process.

# CRISP-DM



Figure 1: Crisp-DM model
https://statistik-dresden.de/archives/1128

# Business Understanding Phase

Our topic: Soccer statistics about the Bundesliga in Germany

Our goal: Understand the Performance development of a team

# Data Understanding Phase

**Source**

Open Data: http://www.football-data.co.uk/

- Contains data for last 10 seasons of German Bundesliga including current season

- Contains various statistical data

- Each season in a separate file

# Necessary tools:

- Latex → perform analytic report

- Microsoft Excel → convert to CSV-files

- PhpMyAdmin → build new database

- MySQL Workbench → design and build DWH/ DM

# GitHub

- Source code management
- Repository
- Work efficently in our team

# Data Preparation Phase

- Data Preparation
- Cleansing
- Compress CSV in ZIP
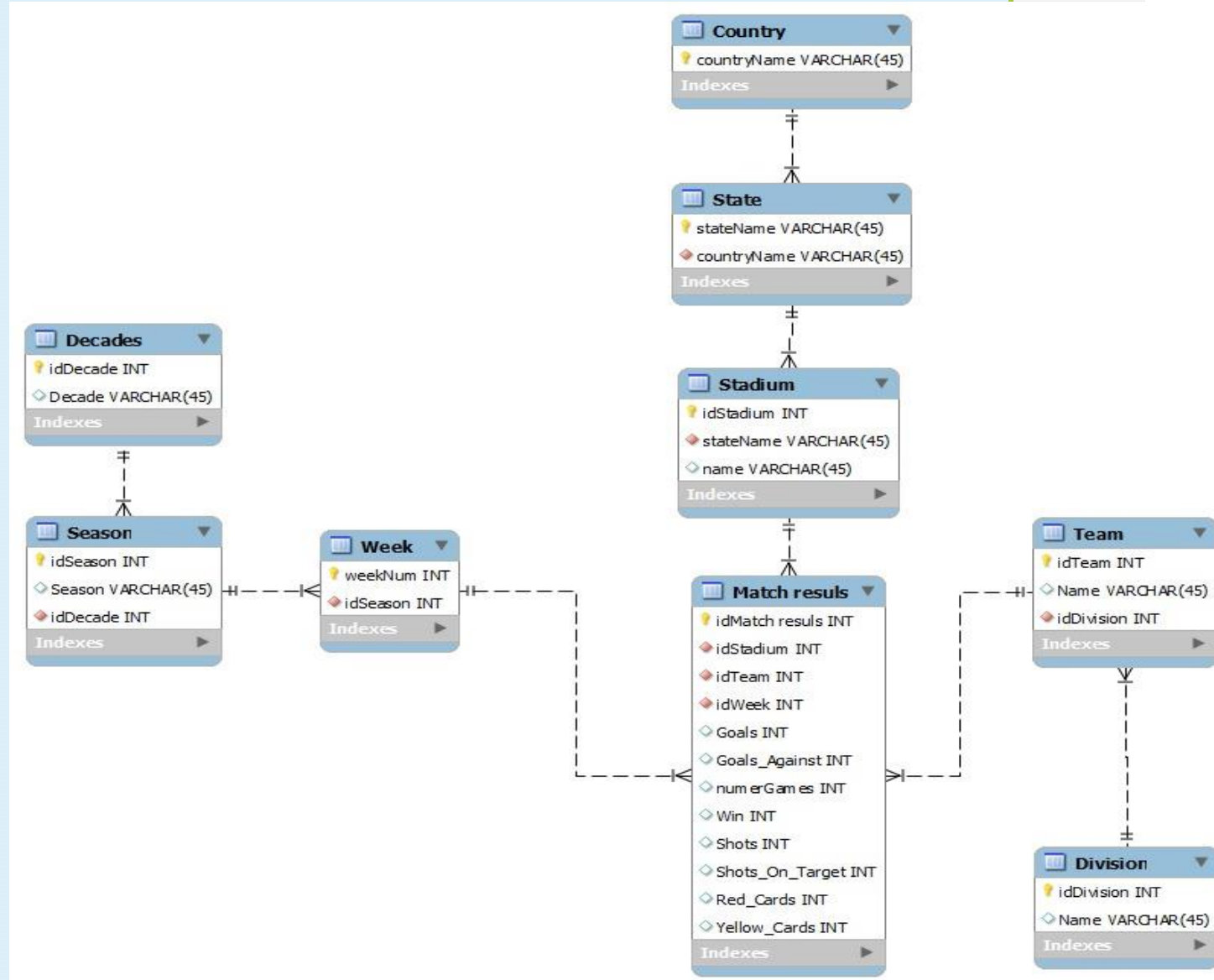- Upload all Seasons (csv.zip)

# Star Schema
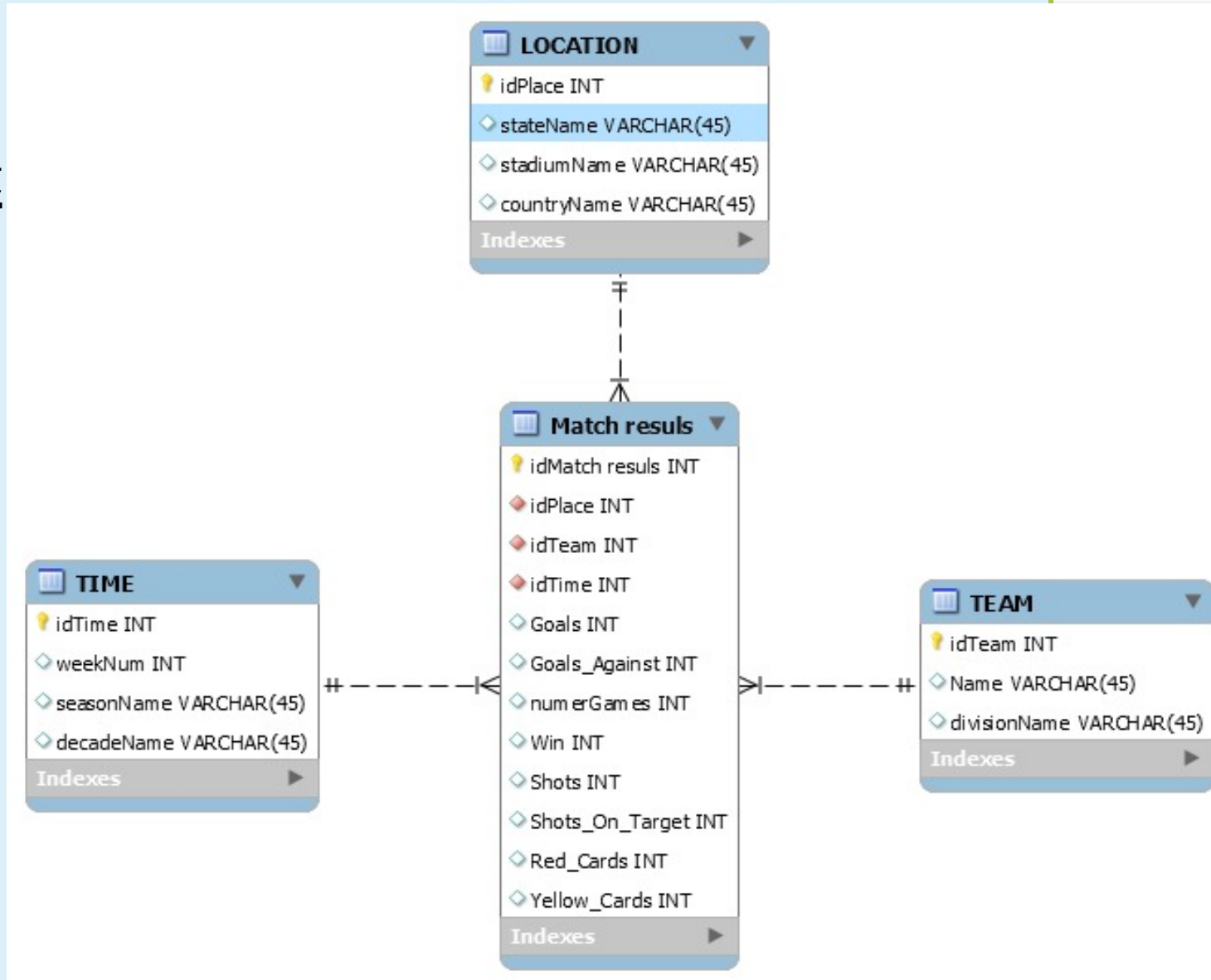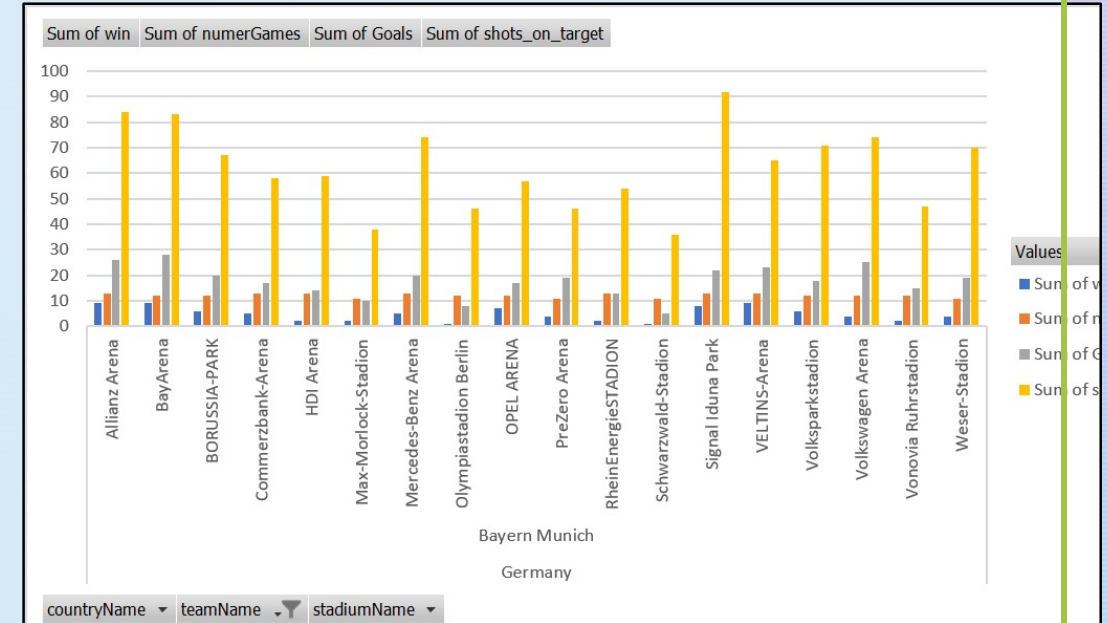


Figure 2: Star Schema

# Data Mart



Figure 3: Data Mart

# Bayern Munich

- Performance of BM
  depending on the stadium



| Row Labels | Sum of win | Sum of numerGames | Sum of Goals | Sum of shots_on_target |
|---|---|---|---|---|
| Germany | 86 | 219 | 319 | 1121 |
| Bayern Munich | 86 | 219 | 319 | 1121 |
| Allianz Arena | 9 | 13 | 26 | 84 |
| BayArena | 9 | 12 | 28 | 83 |
| BORUSSIA-PARK | 6 | 12 | 20 | 67 |
| Commerzbank-Arena | 5 | 13 | 17 | 58 |
| HDI Arena | 2 | 13 | 14 | 59 |
| Max-Morlock-Stadion | 2 | 11 | 10 | 38 |
| Mercedes-Benz Arena | 5 | 13 | 20 | 74 |
| Olympiastadion Berlin | 1 | 12 | 8 | 46 |
| OPEL ARENA | 7 | 12 | 17 | 57 |
| PreZero Arena | 4 | 11 | 19 | 46 |
| RheinEnergieSTADION | 2 | 13 | 13 | 54 |
| Schwarzwald-Stadion | 1 | 11 | 5 | 36 |
| Signal Iduna Park | 8 | 13 | 22 | 92 |
| VELTINS-Arena | 9 | 13 | 23 | 65 |
| Volksparkstadion | 6 | 12 | 18 | 71 |
| Volkswagen Arena | 4 | 12 | 25 | 74 |

# Conclusion

- Raw data is used to create a Data Warehouse
- Check the data quality and reduced to the relevance data
- Easy and compact information about the football statics
- Good a new experience for the group

Thank you...