

Data Quality Assessment:

Dataset: NFL Scores and Betting Data (nfl_teams.csv and spreadspoke_scores.csv)

Source: Kaggle.com, compiled from ESPN, NFL.com, Pro-Football-Reference, NOAA, NFLweather.com, and other betting archives.

1. Data Quality Assessment and Improvement

Accuracy:

Data is pulled from reputable sources and cross-referenced (e.g., betting lines from multiple websites; weather from NOAA and NFLweather.com).

Completeness:

The dataset includes NFL scores from 1966 and betting odds since 1979. However, team_division is missing in some nfl_teams.csv rows.

Consistency:

Several fields require standardization:

- schedule_date is a string, not a date.
- schedule_week is 96% numeric, but 4% are labels (e.g., “Superbowl,” “Wildcard”).
- over_under_line includes 0.4% invalid values (e.g., space characters).
- weather_detail contains inconsistent formatting; 0.8% use pipe delimiters (e.g., "Rain | Fog").

Timeliness:

Data is updated weekly, with a maximum potential 7-day lag post-game.

Uniqueness:

No duplicate rows were found.

Summary:

The dataset is accurate and timely but has minor completeness gaps and several consistency issues involving date, weather, and number formatting.

Transformations / Cleaning Applied:

- In nfl_teams.csv, missing team division values were labeled “not applicable” for now.
 - over_under_line blank spaces were cleaned and converted to numeric.
 - weather_detail was split by delimiter into separate boolean columns (e.g., Rain, Fog).
-

2. Database Schema Design (SQL DDL)

While all of the completeness issues were resolvable by the actions taken in stage 1, to resolve lingering consistency issues, I designed a star schema during ETL, splitting the data into one fact table (Scores) and three dimension tables: Teams, Stadium, and Weather. See fig 1.

- Teams and Stadia are Type 2 dimensions (include is_current, valid_from, valid_to) to handle changing values.
- Weather is a static reference table, as weather types are fixed. New types may be added, but existing types do not have any variables that might change over time.

Date Handling:

Dates stored as varchar(8) in ISO 112 format to prevent regional conflicts during export (e.g., Excel misreading US/UK date formats). This remediates data consistency issues with the date fields.

Rationale:

This schema improves performance (normalisation, more efficient datatypes), data quality (specifically with regards to consistency, as only one column remains with data consistency issues - *schedule_week*), and maintainability. SQL DDL was written and tested using MS SQL Server. Field types and definitions documented below, with sample values included.

Further Improvements:

To turn this into a proper ETL pipeline, the views created would populate staging tables and then a merge usp used to check for differences between staging and production tables, updating each where changes are identified and updating iscurrent and validfrom/validto for deprecated rows. This would allow post and pre-match changes to the tables to be reflected accurately.

schedule_week could be spun out into a schedule dimension table where dates have a match_type dimension and the Score fact table could link to it via a 1:1 PK/SK relationship to resolve the last consistency issue.

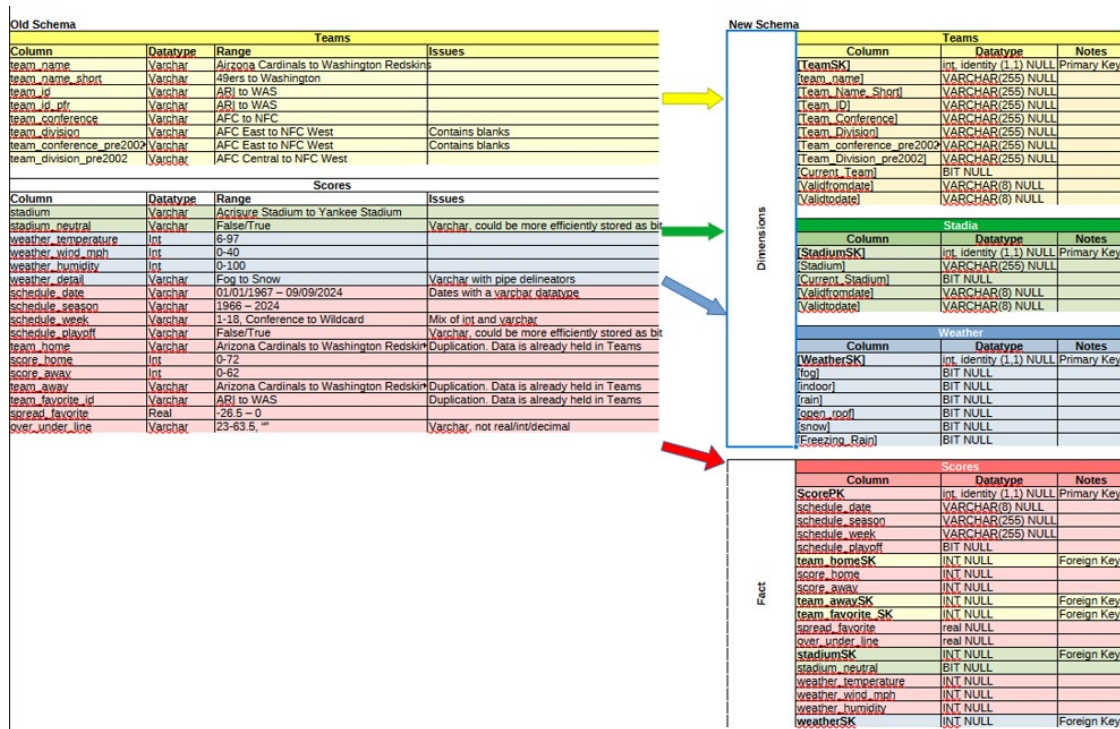


Figure 1: Old Schema to New Schema Overview

Screenshots and SQL Samples:

This submission includes before/after dataset screenshots, schema diagrams, and annotated SQL code showing table creation, data types, and field-level assumptions.

```

1 CREATE TABLE [Teams]
2 (
3 [TeamSK] int identity (1,1),
4 [team_name] VARCHAR(255) NULL,
5 [Team_Name_Short] VARCHAR(255) NULL,
6 [Team_ID] VARCHAR(255) NULL,
7 [Team_Conference] VARCHAR(255) NULL,
8 [Team_Division] VARCHAR(255) NULL,
9 [Team_conference_pre2002] VARCHAR(255) NULL,
10 [Team_Division_pre2002] VARCHAR(255) NULL,
11 [Current_Team] BIT NULL,
12 [Validfromdate] VARCHAR(8) NULL,
13 [Validtodate] VARCHAR(8) NULL
14 );
15
16 CREATE VIEW vpopulate_Teams AS
17 SELECT
18 team_name AS [Team_Name],
19 team_name_short AS [Team_Name_Short],
20 team_id AS [Team_ID],
21 team_conference AS [Team_Conference],
22 team_division AS [Team_Division],
23 team_conference AS [Team_conference_pre2002],
24 team_division_pre2002 AS [Team_Division_pre2002],
25 CASE When team_division = 'Not Applicable' Then 0 Else 1 End AS [Current_Team],
26 19000101 AS [Validfromdate],
27 29991231 AS [Validtodate]
28
29 FROM nfl_teams_amended;
30
31 INSERT INTO Teams
32 SELECT *
33 FROM vpopulate_Teams;
34
35 -- Manual updates for no longer current team names
36
37 UPDATE teams SET validtodate = 31081996 WHERE team_name = 'Houston Oilers';
38 UPDATE teams SET validfromdate = 01091996, validtodate = 31081998 WHERE team_name = 'Tennessee Oilers';
39 UPDATE teams SET validfromdate = 01091998 WHERE team_name = 'Tennessee Titans';
40
41 -- Boston / New England Patriots
42 UPDATE teams SET validtodate = 31081970 WHERE team_name = 'Boston Patriots';
43 UPDATE teams SET validfromdate = 01091971 WHERE team_name = 'New England Patriots';
44
45 -- St. Louis / Phoenix / Arizona Cardinals
46 UPDATE teams SET validtodate = 31081987 WHERE team_name = 'St. Louis Cardinals';
47 UPDATE teams SET validfromdate = 01091988, validtodate = 31081993 WHERE team_name = 'Phoenix Cardinals';
48 UPDATE teams SET validfromdate = 01091994 WHERE team_name = 'Arizona Cardinals';
49
50 -- Baltimore / Indianapolis Colts
51 UPDATE teams SET validtodate = 31081983 WHERE team_name = 'Baltimore Colts';
52 UPDATE teams SET validfromdate = 01091984 WHERE team_name = 'Indianapolis Colts';
53
54 -- Oakland / Los Angeles / Las Vegas Raiders
55 UPDATE teams SET validtodate = 31081981 WHERE team_name = 'Oakland Raiders';
56 UPDATE teams SET validfromdate = 01091982, validtodate = 31081994 WHERE team_name = 'Los Angeles Raiders';
57 UPDATE teams SET validfromdate = 01091995, validtodate = 31082019 WHERE team_name = 'Oakland Raiders';
58 UPDATE teams SET validfromdate = 01092020 WHERE team_name = 'Las Vegas Raiders';
59

```

Figure 2: Teams Table and View Create and populate code

```

70 -- Create Stadia table
71 CREATE Table Stadia
72 (
73 StadiumSK int identity (1,1),
74 Stadium Varchar(255) null,
75 Current_stadium bit null,
76 validfromdate varchar(8),
77 validtodate varchar(8),
78 );
79
80 -- create view to populate it
81 CREATE VIEW vpopulate_Stadia
82 AS
83 SELECT DISTINCT stadium,
84 1 AS current_stadium,
85 19000101 AS validfromdate,
86 29991231 AS validtoDate
87 FROM spreadspoke_scores_amended;
88
89 -- Populate it
90 INSERT INTO Stadia
91 SELECT * FROM vpopulate_Stadia;
92
93
94 -- create a weather reference table
95 CREATE TABLE [Weather]
96 (
97 [WeatherSK] int identity (1,1),
98 [fog] bit null,
99 [indoor] bit null,
100 [rain] bit null,
101 [open_roof] bit null,
102 [snow] bit NULL,
103 [freezing_rain] bit NULL
104 )
105
106 -- create view to populate it with every combination of weather
107 CREATE VIEW vpopulate_Weather
108 AS
109 SELECT a.bit AS fog,
110        b.bit AS indoor,
111        c.bit AS rain,
112        d.bit AS open_roof,
113        e.bit AS snow,
114        f.bit AS freezing_rain
115 FROM (SELECT 0 AS bit UNION ALL SELECT 1) a
116 CROSS JOIN (SELECT 0 AS bit UNION ALL SELECT 1) b
117 CROSS JOIN (SELECT 0 AS bit UNION ALL SELECT 1) c
118 CROSS JOIN (SELECT 0 AS bit UNION ALL SELECT 1) d
119 CROSS JOIN (SELECT 0 AS bit UNION ALL SELECT 1) e
120 CROSS JOIN (SELECT 0 AS bit UNION ALL SELECT 1) f
121
122 -- Populate it
123 INSERT INTO Weather
124 SELECT * FROM vpopulate_Weather
125

```

Figure 3: Weather and Stadia Table and View SQL commands (including a cartesian join for weather)

```

129 CREATE TABLE [Scores]
130 (
131 [ScorePK] int identity (1,1),
132 schedule_date varchar(8) null,
133 schedule_season varchar(4) null,
134 schedule_week Varchar(255) null,
135 schedule_playoff bit null,
136 team_homeSK int null,
137 score_home int null,
138 score_away int null,
139 team_awaySK int null,
140 team_favoriteSK int null,
141 spread_favorite real null,
142 over_under_line real null,
143 stadiumsk int null,
144 stadium_neutral bit null,
145 weather_temperature int null,
146 weather_wind_mph int null,
147 weather_humidity int null,
148 weatherSK int null)
149
150 CREATE VIEW vpopulate_Scores AS
151
152 SELECT
153 Convert (Varchar(8), Cast (s.schedule_date AS Date),112) AS schedule_date,
154 s.schedule_season AS schedule_season,
155 s.schedule_week AS schedule_week,
156 CASE WHEN s.schedule_playoff = 'FALSE' then 0 ELSE 1 END AS schedule_playoff,
157 ht.TeamSK AS team_homesk,
158 s.score_home AS score_home,
159 s.score_away AS score_away,
160 at.TeamSK AS team_awaySK,
161 ft.TeamSK AS team_favoritesk,
162 s.spread_favorite AS spread_favoriteSK,
163 s.over_under_line AS over_under_line,
164 st.StadiumSK AS stadiumSK,
165 CASE WHEN s.stadium_neutral = 'FALSE' then 0 ELSE 1 END AS stadium_neutral,
166 s.weather_temperature AS weather_temperature,
167 s.weather_wind_mph AS weather_wind_mph,
168 s.weather_humidity AS weather_humidity,
169 w.WeatherSK AS weathersk
170
171 FROM spreadspoke_scores_amended S
172 LEFT JOIN teams ht ON s.team_home = ht.team_name
173 AND CONVERT(Varchar(8), Cast (s.schedule_date AS Date),112) >= ht.Validfromdate
174 AND convert(Varchar(8), Cast (s.schedule_date AS Date),112) < ht.Validtodate
175 LEFT JOIN teams at ON s.team_away = at.team_name
176 AND CONVERT(Varchar(8), Cast (s.schedule_date AS Date),112) >= at.Validfromdate
177 AND convert(Varchar(8), Cast (s.schedule_date AS Date),112) < at.Validtodate
178 LEFT JOIN teams ft ON s.team_favorite_id = ft.team_id
179 AND CONVERT(Varchar(8), Cast (s.schedule_date AS Date),112) >= ft.Validfromdate
180 AND convert(Varchar(8), Cast (s.schedule_date AS Date),112) < ft.Validtodate
181 LEFT JOIN Weather w ON ISNULL(s.weather_Fog, 0) = ISNULL(w.fog, 0)
182 AND ISNULL(s.weather_Indoor, 0) = ISNULL(w.indoor, 0)
183 AND ISNULL(s.weather_Rain, 0) = ISNULL(w.rain, 0)
184 AND ISNULL(s.open_Roof, 0) = ISNULL(w.open_roof, 0)
185 AND ISNULL(s.weather_Snow, 0) = ISNULL(w.snow, 0)
186 AND ISNULL(s.weather_FreezingRain, 0) = ISNULL(w.freezing_rain, 0)
187 Left Join Stadia st ON s.stadium = st.stadium
188 AND CONVERT(Varchar(8), CAST(s.schedule_date AS date),112) between st.validfromdate and st.validtodate

```

Figure 4: Scores Fact Table and View Create

	schedule_...	schedule_season	sche...	schedul...	team_home	score...	score...	team_away	team...	sprea...	over_...	stadi...	stadi...	weat...	weat...	weat...	weat...
1/5/2025	2024	18	FALSE		Pittsburgh Steelers	17	19	Cincinnati Bengals	CIN	-2.5	48.5	Acrisur...	FALSE	NULL	NULL	NULL	
1/5/2025	2024	18	FALSE		Tampa Bay Buccaneers	27	19	New Orleans Saints	TB	-14.5	44.5	Raymo...	FALSE	NULL	NULL	NULL	
1/5/2025	2024	18	FALSE		Tennessee Titans	14	23	Houston Texans	TEN	-2.5	36.5	Nissan...	FALSE	NULL	NULL	NULL	
1/11/2025	2024		Wildcard	TRUE	Houston Texans	32	12	Los Angeles Chargers	LAC	-2.5	42	NRG S...	FALSE	72	0	NULL	indoor
1/12/2025	2024		Wildcard	TRUE	Baltimore Ravens	28	14	Pittsburgh Steelers	BAL	-9	44.5	M&T B...	FALSE	32	12	54	
1/12/2025	2024		Wildcard	TRUE	Buffalo Bills	31	7	Denver Broncos	BUF	-7.5	49	Highm...	FALSE	32	10	76	
1/12/2025	2024		Wildcard	TRUE	Philadelphia Eagles	22	10	Green Bay Packers	PHI	-5.5	46	Lincoln...	FALSE	40	6	47	
1/12/2025	2024		Wildcard	TRUE	Tampa Bay Buccaneers	20	23	Washington Commanders	TB	-3	51.5	Raymo...	FALSE	58	2	57	
1/13/2025	2024		Wildcard	TRUE	Los Angeles Rams	27	9	Minnesota Vikings	MIN	-2.5	48	State F...	TRUE	72	0	NULL	indoor
1/18/2025	2024		Division	TRUE	Detroit Lions	31	45	Washington Commanders	DET	-8	55	Ford Fl...	FALSE	72	0	NULL	indoor
1/18/2025	2024		Division	TRUE	Kansas City Chiefs	23	14	Houston Texans	KC	-9.5	41.5	GEHA ...	FALSE	23	17	43	
1/19/2025	2024		Division	TRUE	Buffalo Bills	27	25	Baltimore Ravens	BAL	-1.5	52	Highm...	FALSE	19	4	84	snow
1/19/2025	2024		Division	TRUE	Philadelphia Eagles	28	22	Los Angeles Rams	PHI	-7	43.5	Lincoln...	FALSE	34	8	67	snow

Figure 5: Old Scores Table

	ScorePK	schedule_...	schedule_...	schedule_...	schedule_...	team_hom...	score_home	score_away	team_awa...	team_favo...	spread_fa...	over_unde...	stadiumsk	stadium_n...	weather_t...	weather_w...	weather_h...	weatherSK
73	20250112	2024	Wildcard	1	35	28	14	38	35	-9	44.5	21	0	32	12	54	1	
74	20250112	2024	Wildcard	1	4	31	7	5	4	-7.5	49	14	0	32	10	76	1	
75	20250112	2024	Wildcard	1	18	22	10	26	18	-5.5	46	18	0	40	6	47	1	
76	20250112	2024	Wildcard	1	7	20	23	42	7	-3	51.5	31	0	58	2	57	1	
77	20250113	2024	Wildcard	1	33	27	9	41	41	-2.5	48	36	1	72	0	NULL	2	
78	20250118	2024	Division	1	23	31	45	42	23	-8	55	10	0	72	0	NULL	2	
79	20250118	2024	Division	1	13	23	14	39	13	-9.5	41.5	12	0	23	17	43	1	
80	20250119	2024	Division	1	4	27	25	35	35	-1.5	52	14	0	19	4	84	9	
81	20250119	2024	Division	1	18	28	22	33	18	-7	43.5	18	0	34	8	67	9	
82	20250126	2024	Conference	1	18	55	23	42	18	-6	47	18	0	40	11	45	1	
83	20250126	2024	Conference	1	13	32	29	4	13	-1	49.5	12	0	31	8	45	1	
84	20250209	2024	Superbowl	1	18	40	22	13	13	-1	48.5	6	1	72	0	0	2	

Figure 6: New Scores Table