

TRATAMIENTO DE DATOS OPTA

PFM - MASTER EN BIG DATA DEPORTIVO



ÍNDICE

1. Planteamiento del problema
2. Solución propuesta
3. Fuentes de datos
4. Arquitectura del proyecto
5. Limpieza, Ingesta y Procesamiento de Datos
6. Visualización de Resultados
7. Conclusiones y líneas futuras

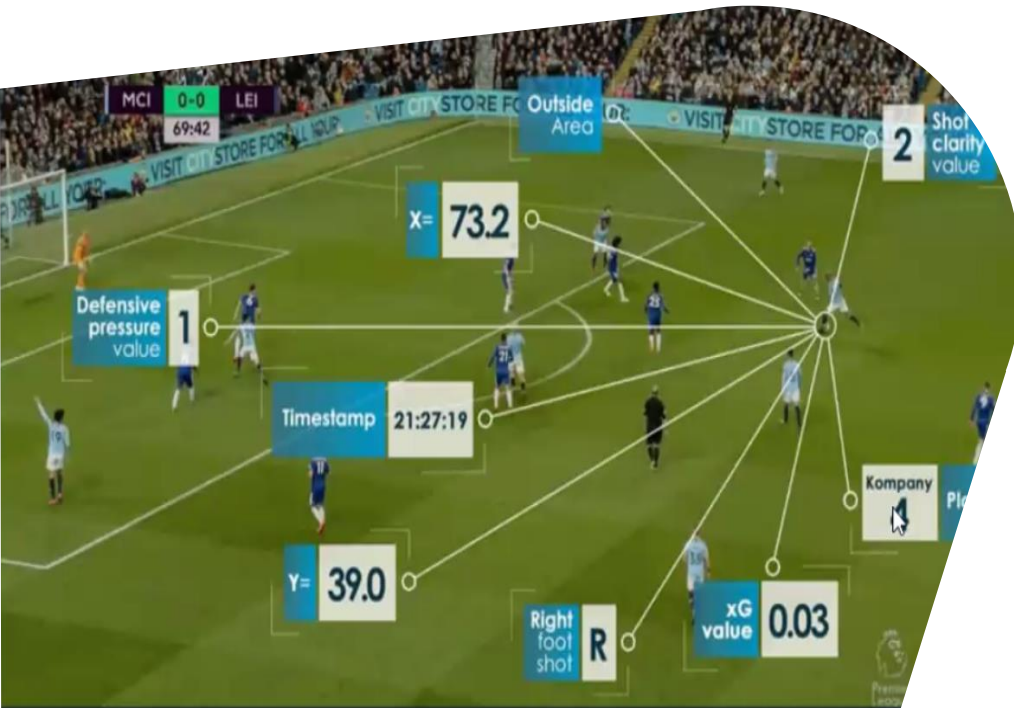
PLANTEAMIENTO DEL PROBLEMA

CONTEXTO Y MOTIVACIONES

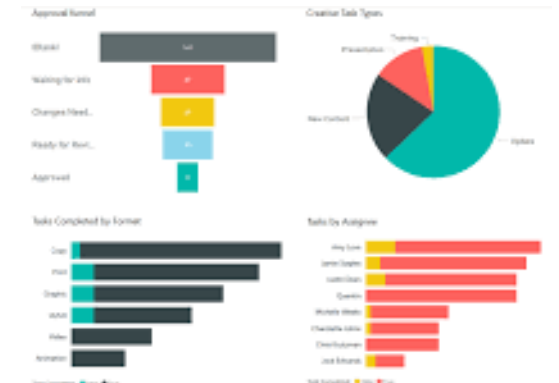
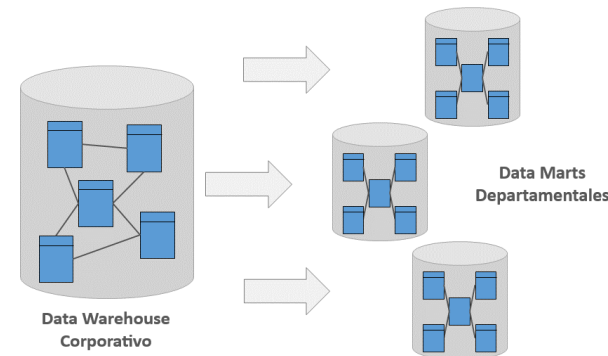
- _ OPTA posee un amplio conjunto de **feeds de datos** de fútbol
- _ Sus clientes tienen acceso a estos datos en **formato XML**



- _ Se desea almacenar toda esta información de manera estructurada en una base de datos (**data warehouse/data mart**, según su alcance) que se pueda explotar mediante **visualizaciones analíticas**



Cuando Kompany golpeó el balón para marcar contra el Leicester, inmovilizó ese momento a través de distintas variables en tiempo real

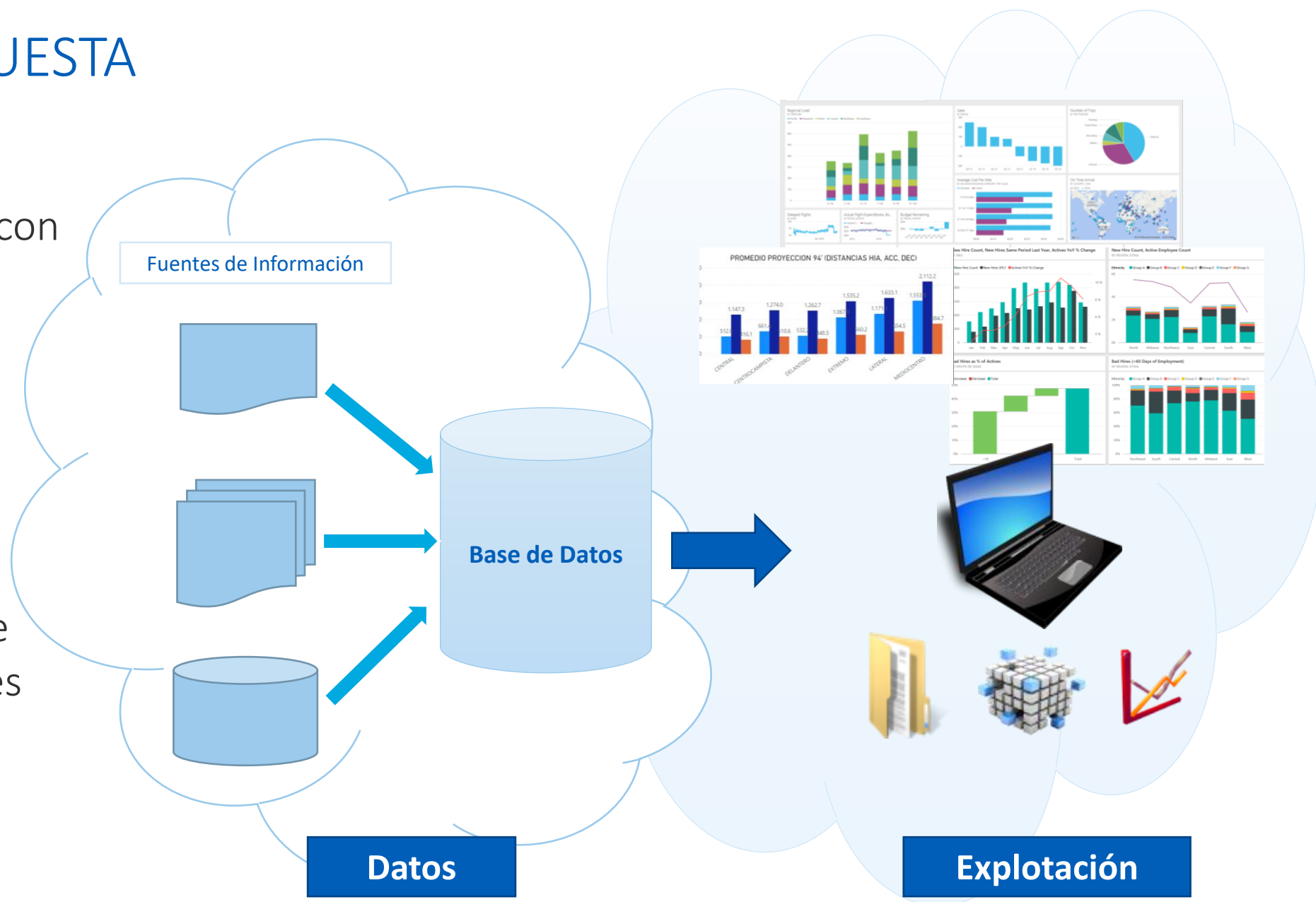


SOLUCIÓN PROPUESTA

_ Nueva base de datos con información de OPTA

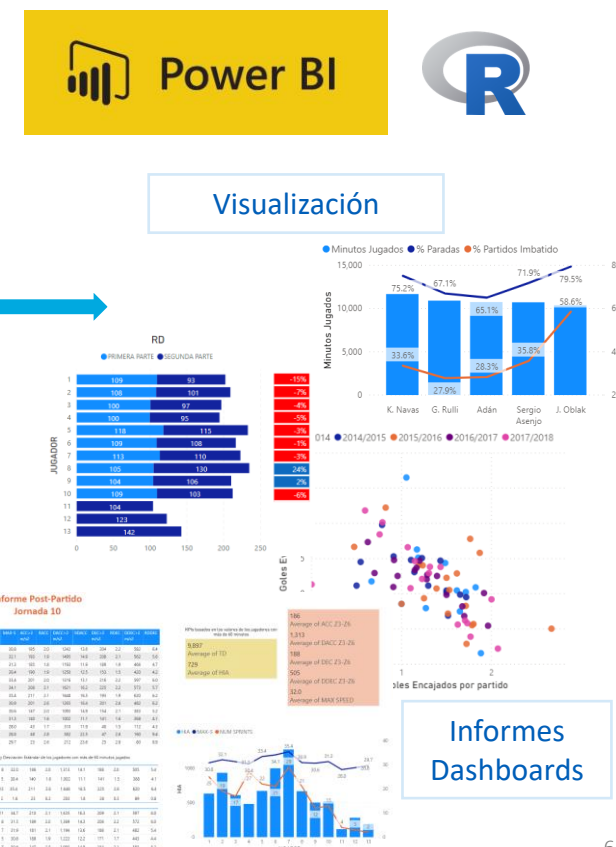
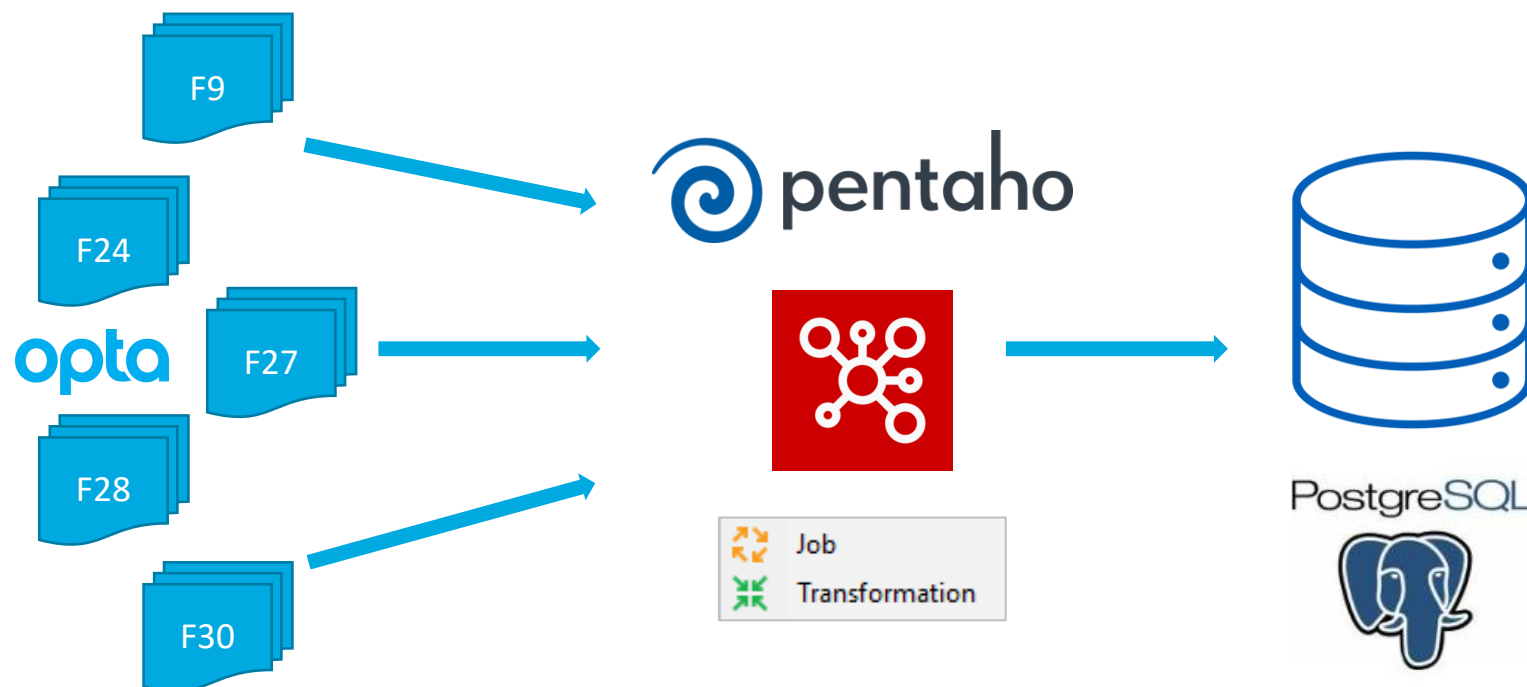
_ Procesos de carga

_ Visualización de la información mediante dashboards o informes



ARQUITECTURA DEL PROYECTO

TECNOLOGIAS UTILIZADAS



(*) Se ha optado por PostgreSQL como gestor de base de datos pero se puede adaptar a cualquier otro.
Pentaho Data Integration (PDI) ofrece conectores para gran cantidad de proveedores

LIMPIEZA, INGESTA Y PROCESAMIENTO DE DATOS

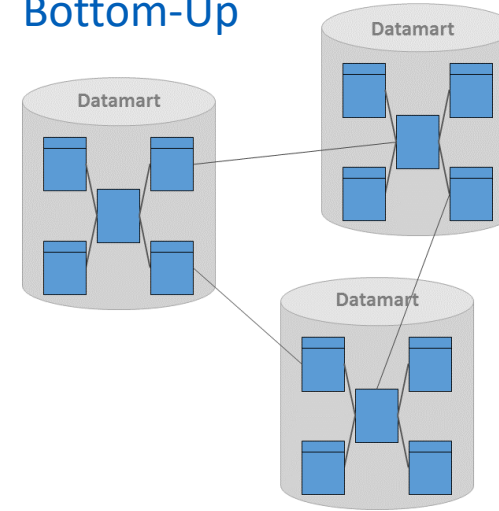
CONCEPTOS TÉCNICOS Y METODOLOGÍAS

_ DATA WAREHOUSE (DW): repositorio de datos corporativos estratégicos, tácticos y operativos. Información de interés generada por la actividad de una empresa u organización

_ DATA MART (DM): subconjunto físico/lógico del DW preparado para la consulta y análisis de la información de un área específico del negocio

_ PROCESOS ETL (Extract, Transform y Load)

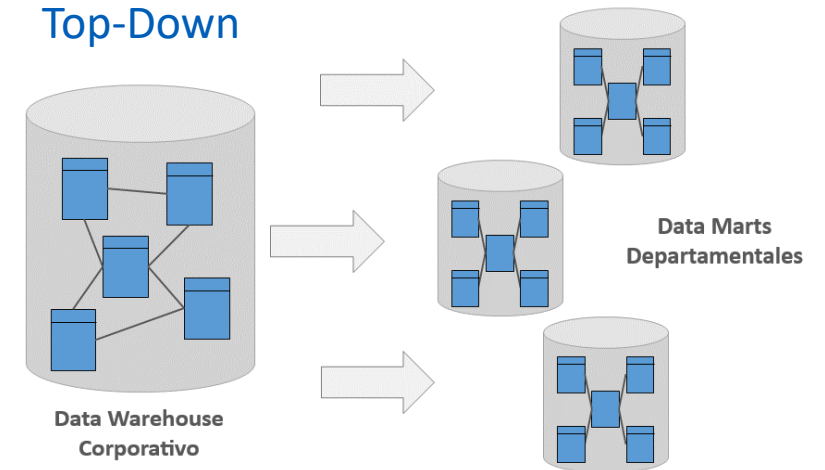
Bottom-Up



Características:

- Orientado a temas
- Integrado
- Variante en el tiempo
- No volátil

Top-Down



LIMPIEZA, INGESTA Y PROCESAMIENTO DE DATOS

CONCEPTOS TÉCNICOS Y METODOLOGÍAS

_MODELO MULTIDIMENSIONAL

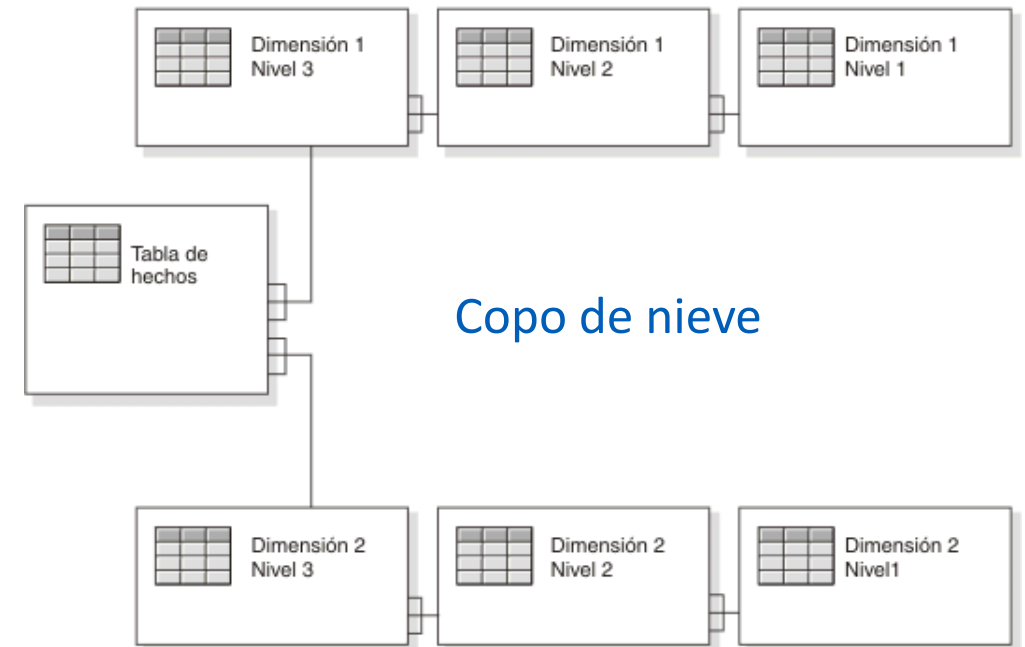
- Hechos
- Dimensiones

_ESQUEMAS DE DISEÑO

- Estrella
- Copo de nieve



Estrella



Copo de nieve

LIMPIEZA, INGESTA Y PROCESAMIENTO DE DATOS



PROCESOS ETL

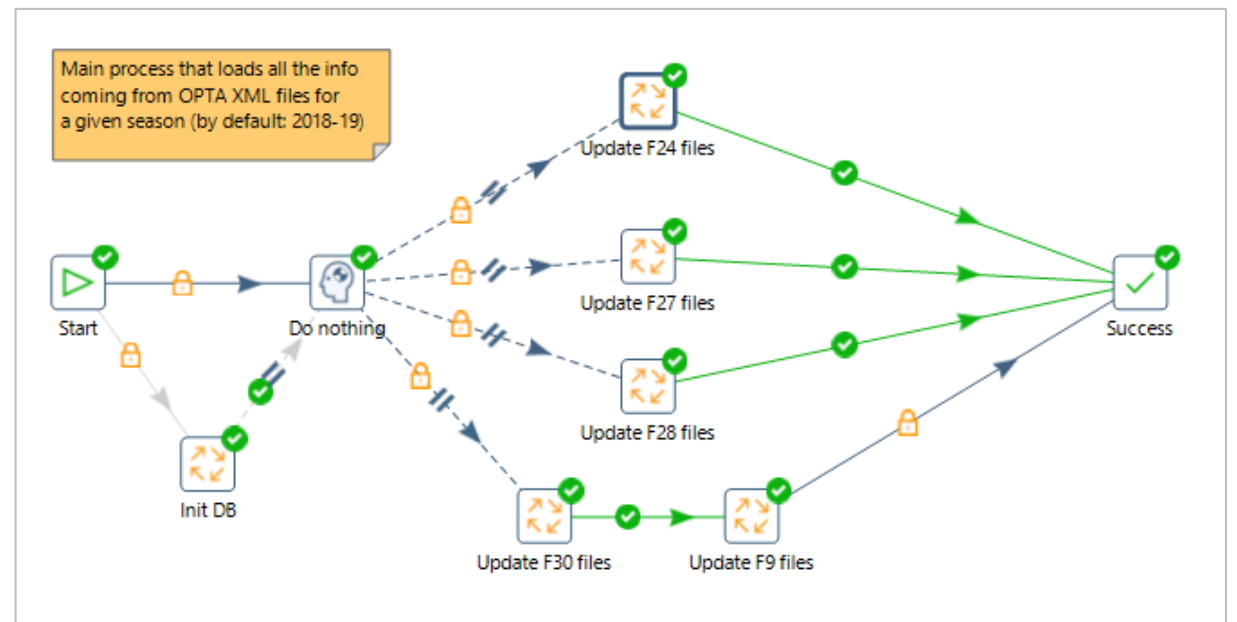
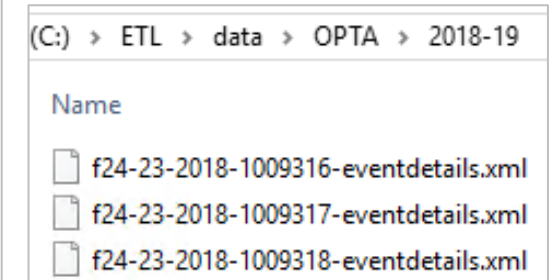
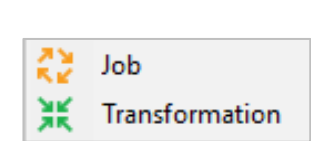
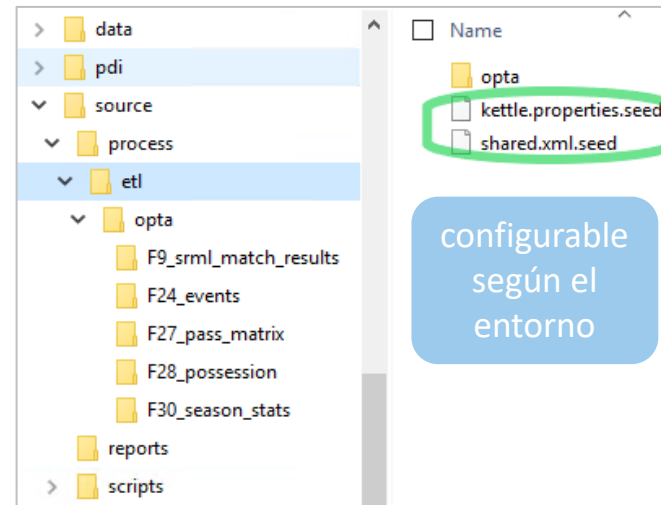
_Pentaho Data Integration (PDI)

- Jobs / Trabajos
- Transformaciones

_Procesos parametrizados con los datos de la temporada (ej. *season = 2018-19*)

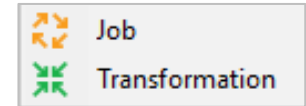
- El nombre del directorio que contiene los XML deberá ajustarse a este parámetro
- Ficheros *kettle.properties* y *shared.xml*

_Proceso/job principal “orquesta” la ejecución del resto de procesos



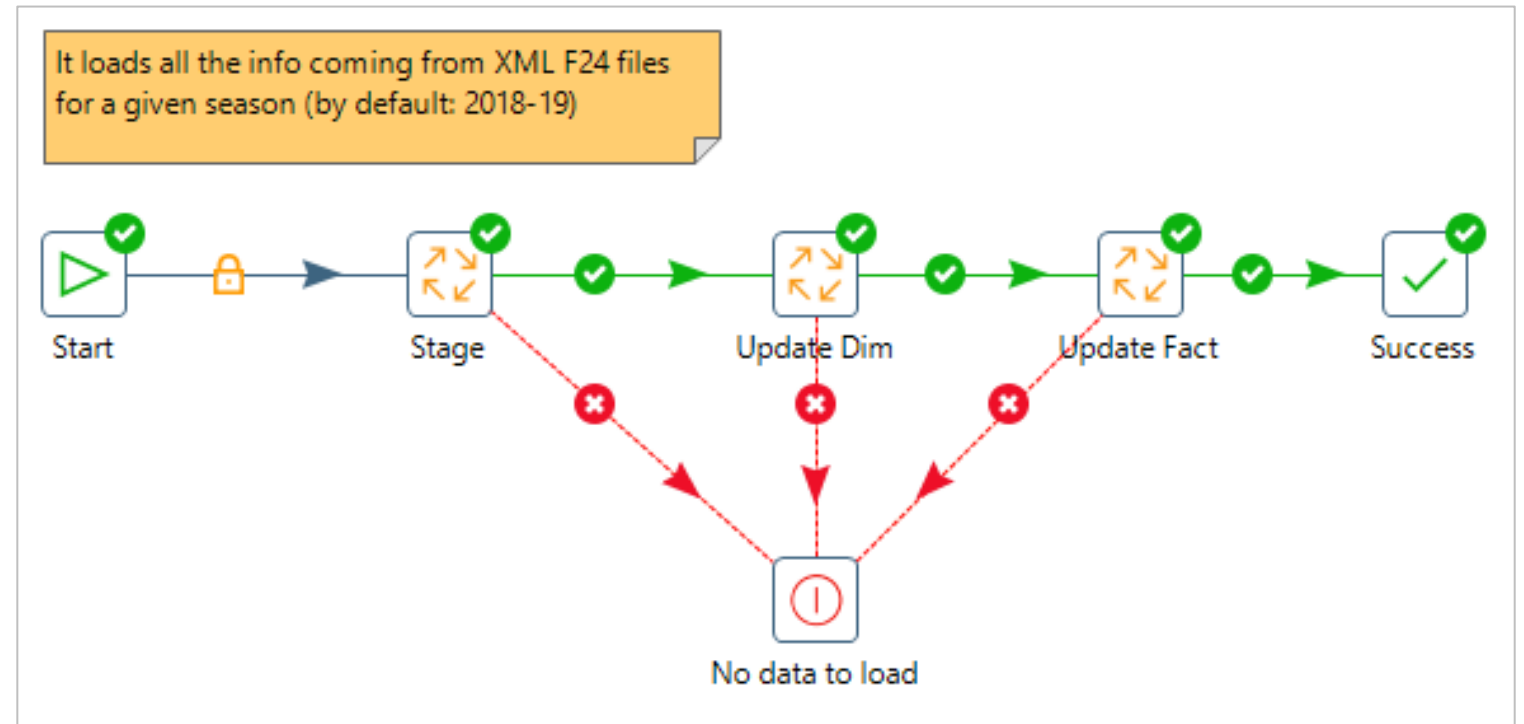
LIMPIEZA, INGESTA Y PROCESAMIENTO DE DATOS

PROCESOS ETL



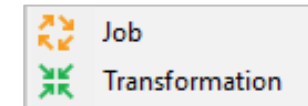
_Metodología común para todos los XML de OPTA con 3 trabajos (.kjb) y varias transformaciones (.ktr):

- Stage area, carga XML en tablas temporales
- Actualización de dimensiones (dim)
- Actualización de tablas de hechos (fact)



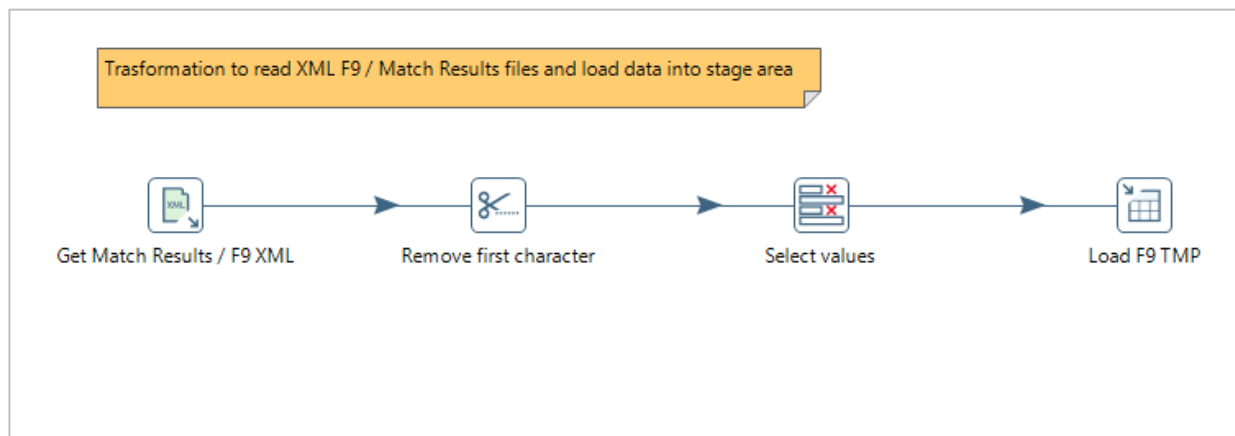
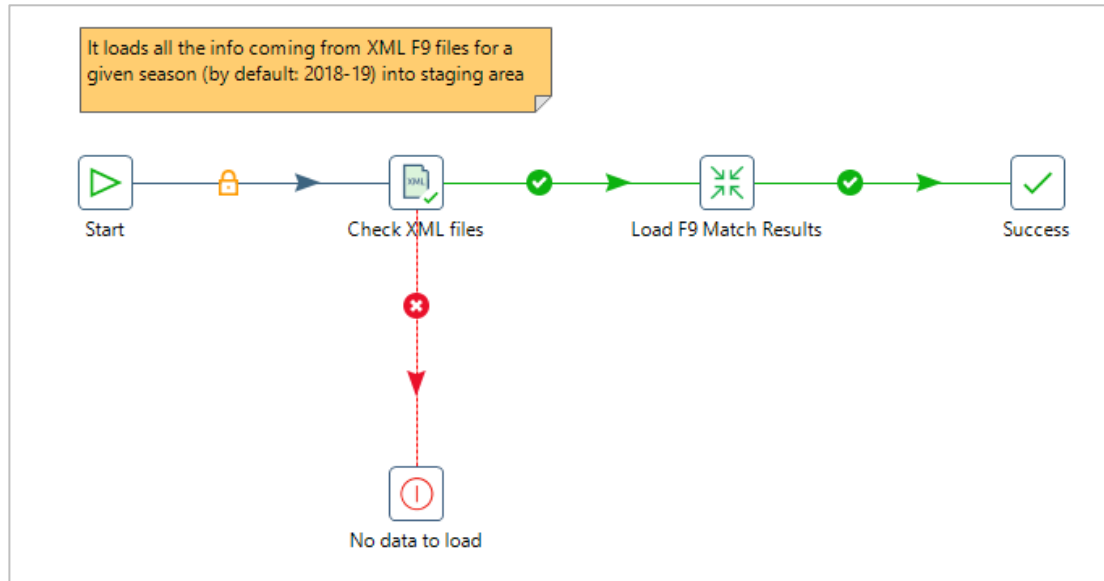
LIMPIEZA, INGESTA Y PROCESAMIENTO DE DATOS

PROCESOS ETL (Stage)



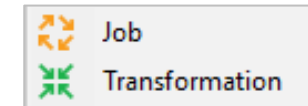
Jobs STAGE

1. Comprobación de la existencia de ficheros XML
 - Error en caso de no haber datos
2. Transformación para cargar el contenido en el data stage de la base de datos
 - Lectura del fichero XML
 - Tratamiento de datos, si procede (creación de nuevas columnas, substring, pivotado...)
 - Inserción en tabla temporal

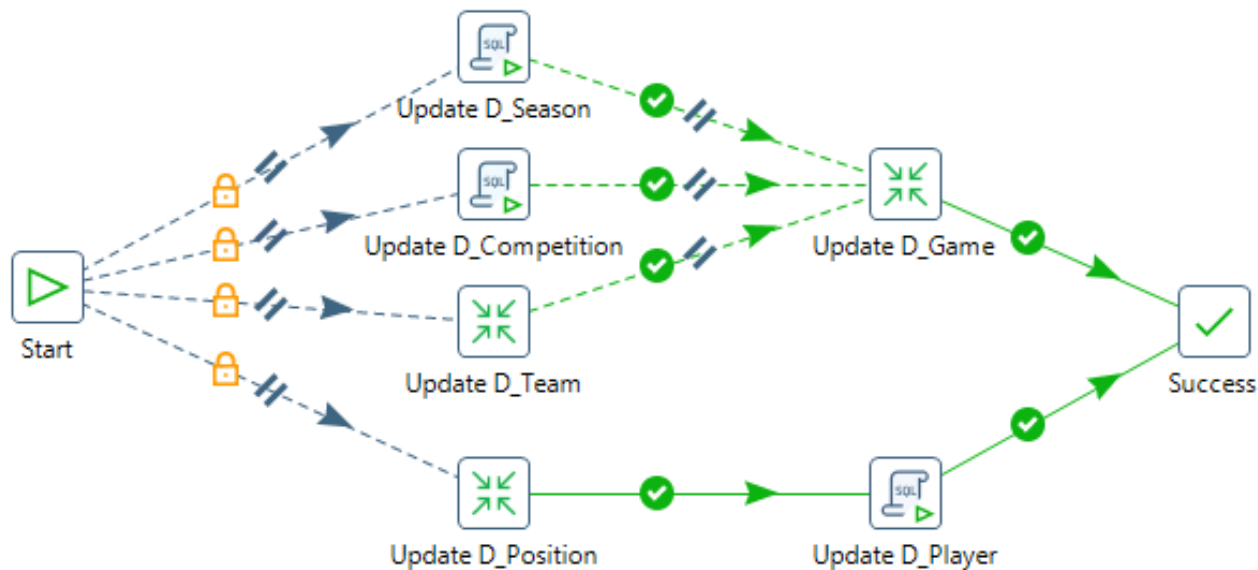


LIMPIEZA, INGESTA Y PROCESAMIENTO DE DATOS

PROCESOS ETL (Dimensiones)



It updates the dimension values according to the info coming from TMP F27 table

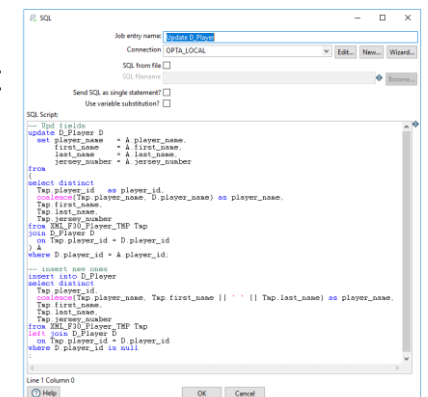


Update DIMs

- _ Actualización de atributos
- _ Inserción de nuevos elementos que aparecen en los ficheros a cargar

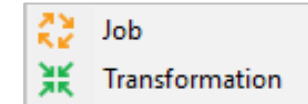
mediante

- Transformaciones específicas ofrecidas por PDI
- Queries SQL ad-hoc



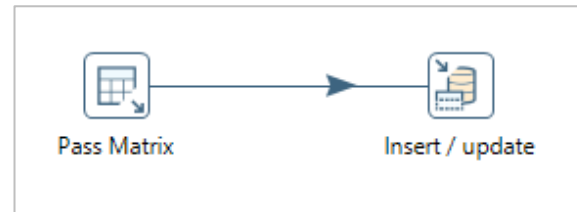
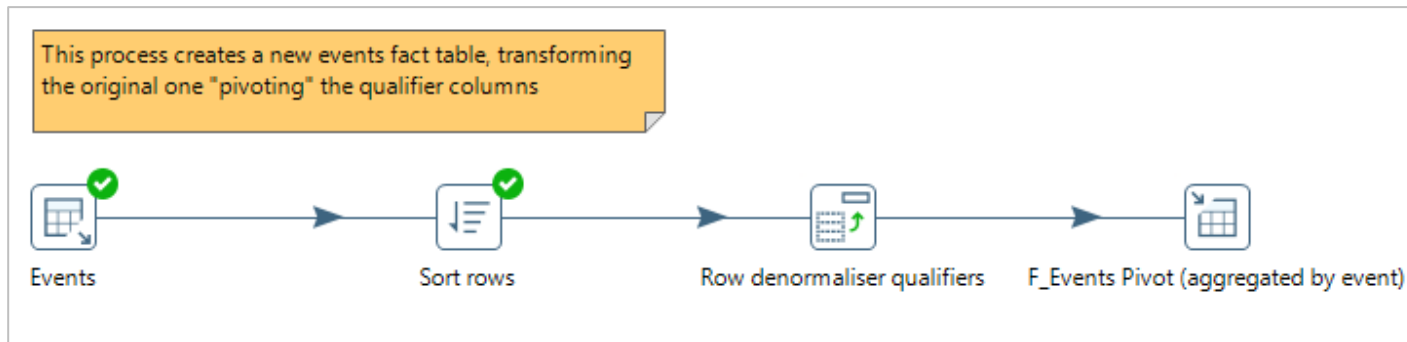
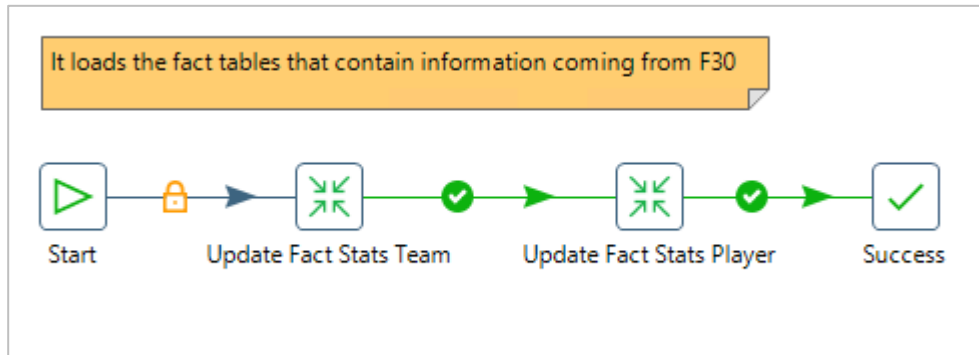
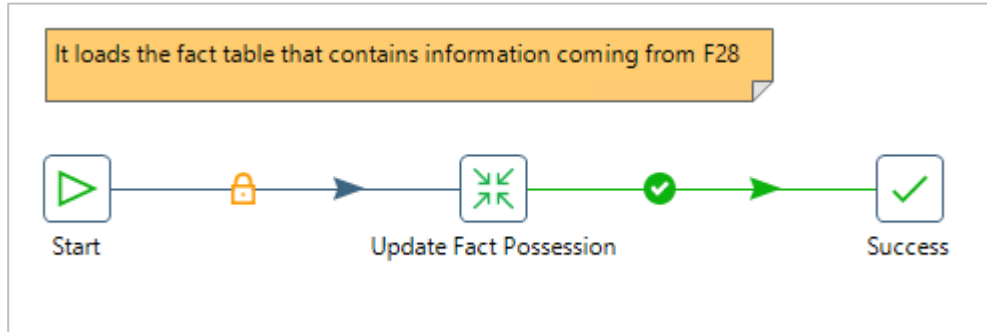
LIMPIEZA, INGESTA Y PROCESAMIENTO DE DATOS

PROCESOS ETL (Hechos)



Update FACTs

Sin borrados, excepto en agregados/f24 pivotado, solo inserción de nuevos datos o actualización de los existentes en caso de recargas



Insert / update

Step name: Insert / update

Connection: OPTA_LOCAL

Target schema: []

Target table: F_Pass_Matrix

Commit size: 10000

Don't perform any updates: ☐

The key(s) to look up the value(s):

#	Table field	Comparator	Stream field1	Stream field2
1	game_id	=	game_id	
2	team_id	=	team_id	
3	passer_player_id	=	passer_player_id	
4	receiver_player_id	=	receiver_player_id	

Update fields:

#	Table field	Stream field	Update
1	game_id	game_id	N
2	team_id	team_id	N
3	passer_player_id	passer_player_id	N
4	receiver_player_id	receiver_player_id	N
5	position_id	position_id	Y
6	x	x	Y
7	y	y	Y
8	cross_success	cross_success	Y
9	pass_success	pass_success	Y

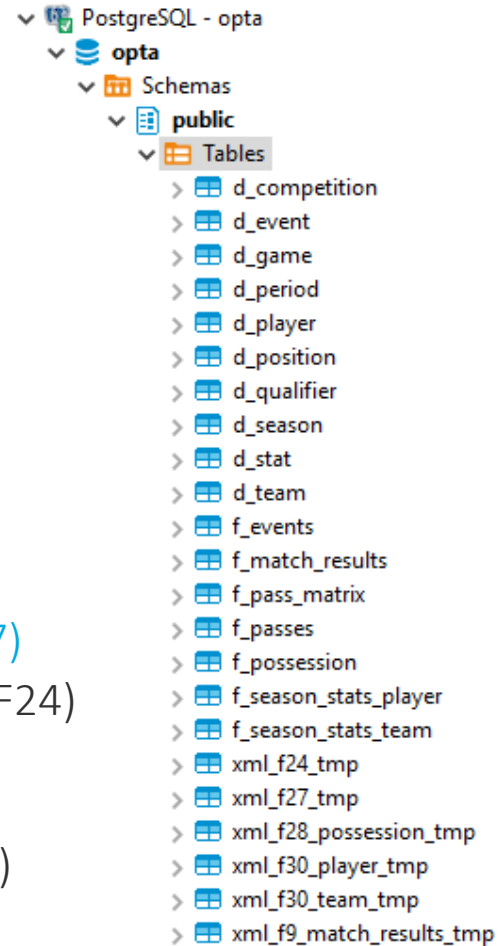
LIMPIEZA, INGESTA Y PROCESAMIENTO DE DATOS

MODELO DE DATOS



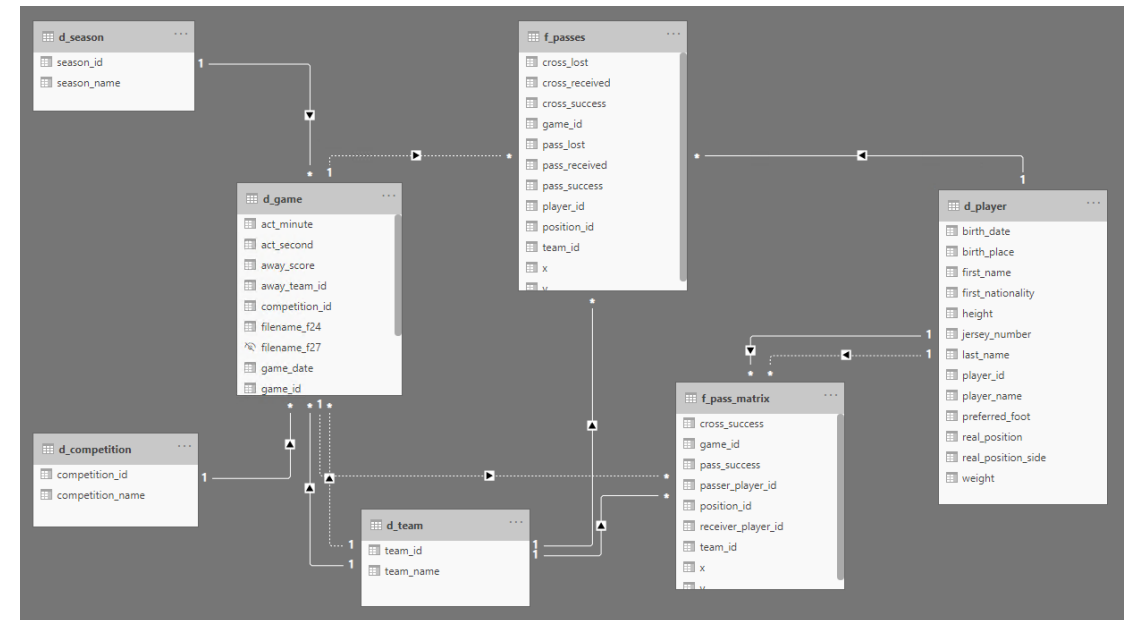
_Tablas de Dimensiones

- D_Competition
- D_Season
- D_Game
- D_Team
- D_Player
- D_Position
- D_Stat
- D_Qualifier



_Tablas de Hechos

- F_Pass y F_Pass_Matrix (F27)
- F_Events y F_Events_Pivot (F24)
- F_Possession (F28)
- F_Season_Stats_Team y F_Season_Stats_Player (F30)
- F_Match_Results (F9)



LIMPIEZA, INGESTA Y PROCESAMIENTO DE DATOS

MODELO DE DATOS

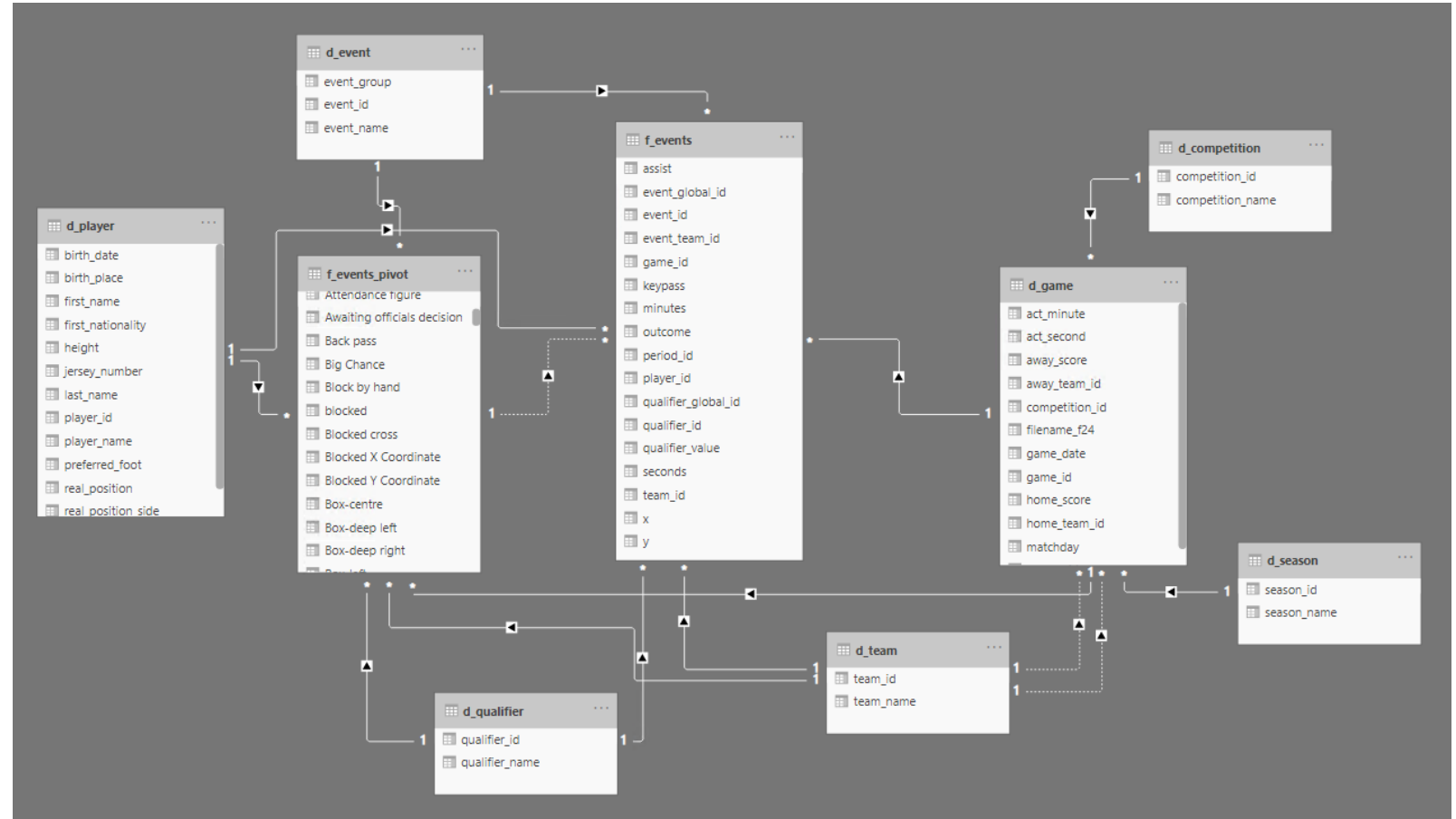


_Tablas de Dimensiones

- D_Competition
- D_Season
- D_Game
- D_Team
- D_Player
- D_Position
- D_Stat
- D_Event
- D_Qualifier

_Tablas de Hechos

- F_Pass y F_Pass_Matrix (F27)
- F_Events y F_Events_Pivot (F24)
- F_Possession (F28)
- F_Season_Stats_Team y F_Season_Stats_Player (F30)
- F_Match_Results (F9)



LIMPIEZA, INGESTA Y PROCESAMIENTO DE DATOS

MODELO DE DATOS

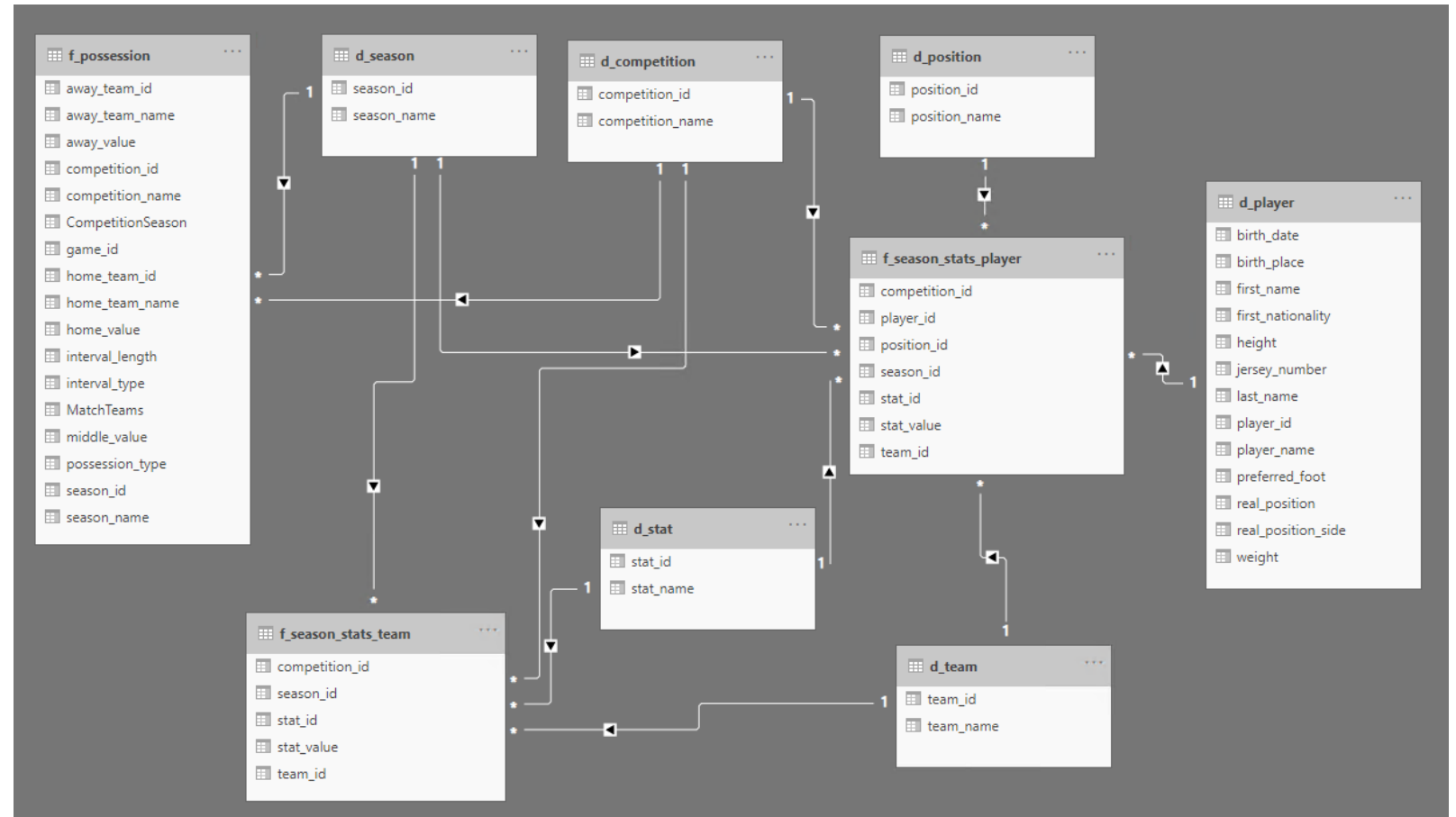


_Tablas de Dimensiones

- D_Competition
- D_Season
- D_Game
- D_Team
- D_Player
- D_Position
- D_Stat
- D_Event
- D_Qualifier

_Tablas de Hechos

- F_Pass y F_Pass_Matrix (F27)
- F_Events y F_Events_Pivot (F24)
- **F_Possession (F28)**
- **F_Season_Stats_Team y F_Season_Stats_Player (F30)**
- F_Match_Results (F9)



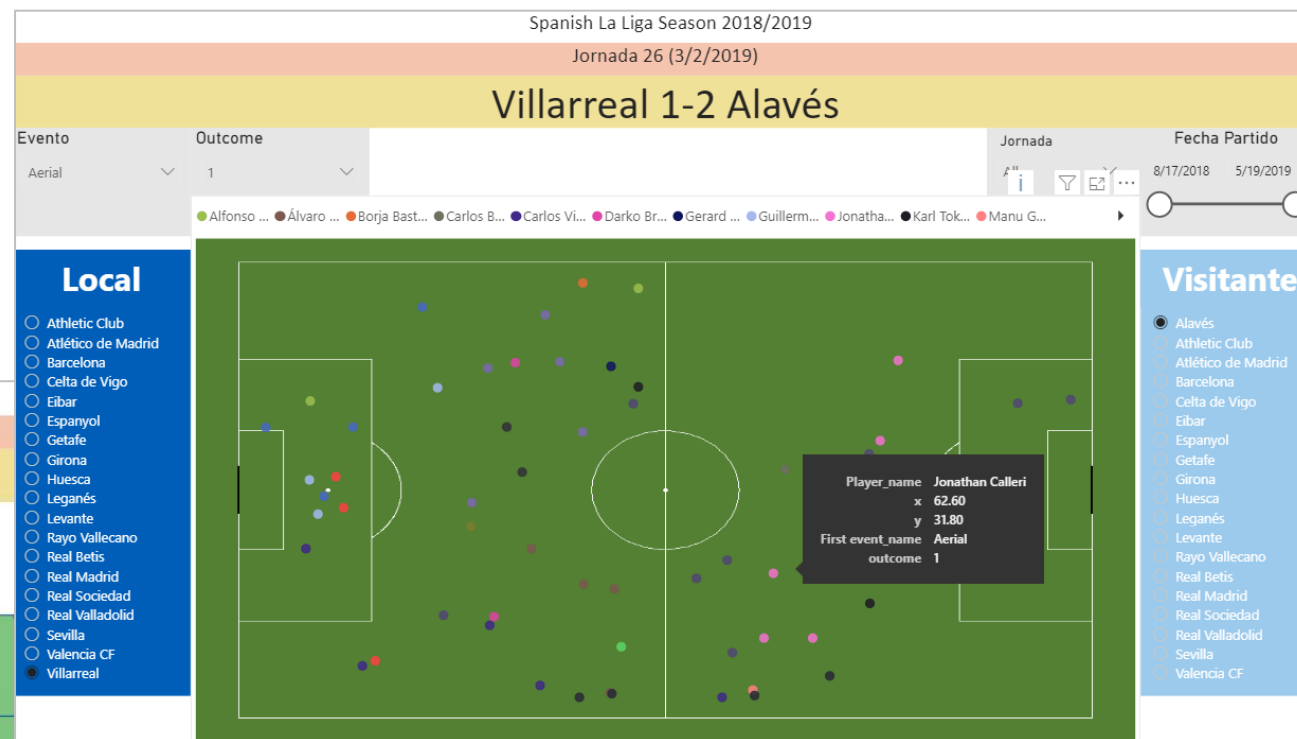
VISUALIZACIÓN DE RESULTADOS

POWER BI (*Events.pbix*)



_ Campograma con eventos

_ Áreas convexas



Área convexa
con eventos seleccionados
(integración de R con Power BI)

VISUALIZACIÓN DE RESULTADOS

POWER BI (Events.pbix)

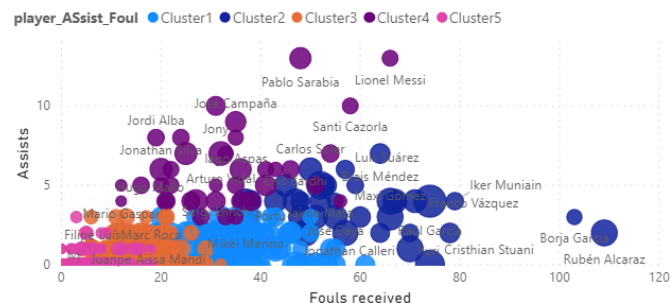


_ Pases

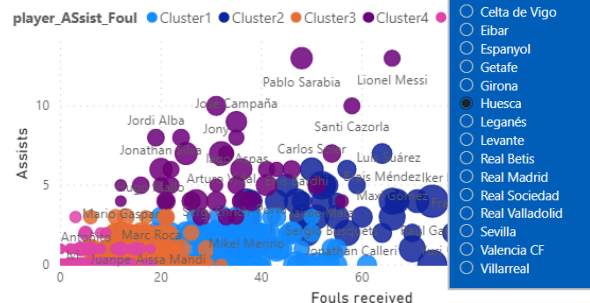
_ Duelos aéreos

_ Cluster de jugadores

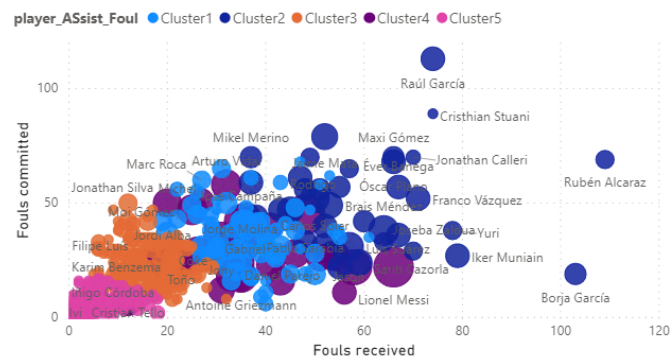
Fouls received, Assists and Fouls committed by player_name and player_ASsist_Foul



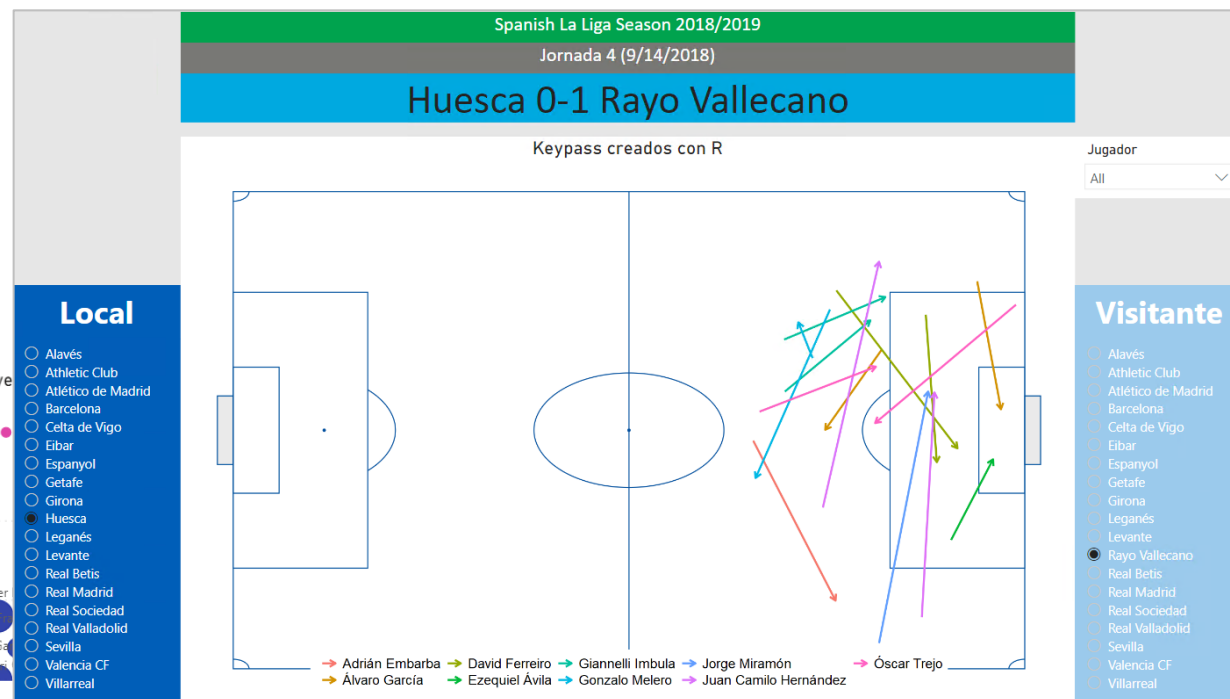
Fouls received, Assists and Fouls committed by player_name and player_ASsist_Foul



Fouls received, Fouls committed and Assists by player_name and player_ASsist_Foul



Player	Assists	Fouls committed	Fouls received	player_ASsist_Foul
Lionel Messi	13	22	66	Cluster4
Pablo Sarabia	13	45	48	Cluster4
Jony	10	34	31	Cluster4
Santi Cazorla	10	22	58	Cluster4
José Campaña	9	42	35	Cluster4
Antoine Griezmann	8	18	35	Cluster4
Jordi Alba	8	25	19	Cluster4
Wissam Ben Yedder	8	26	24	Cluster4
Arturo Vidal	7	58	32	Cluster4
Brais Méndez	7	38	64	Cluster2
Daniel Parejo	7	27	54	Cluster4
Moi Gómez	7	47	25	Cluster4
Sergi Roberto	7	18	33	Cluster4
Total	615	10306	9734	



VISUALIZACIÓN DE RESULTADOS

POWER BI (Stats & Possession.pbix)

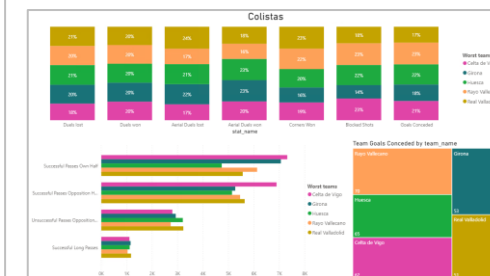
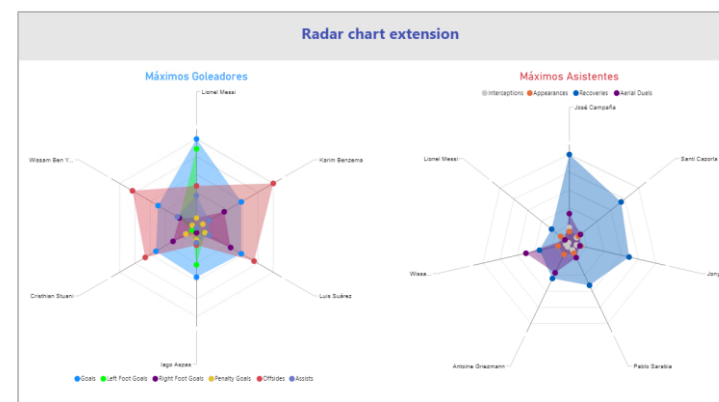
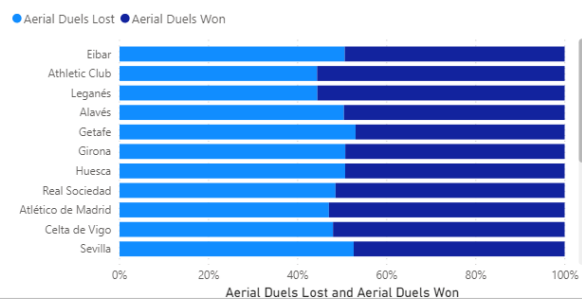
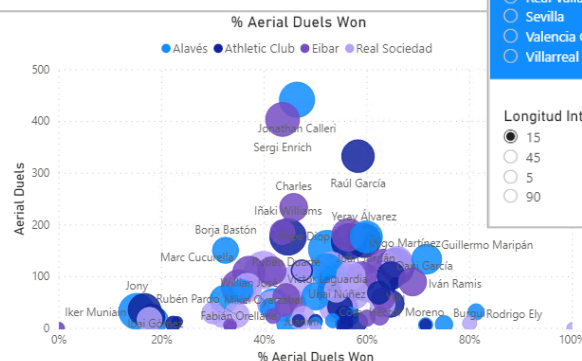
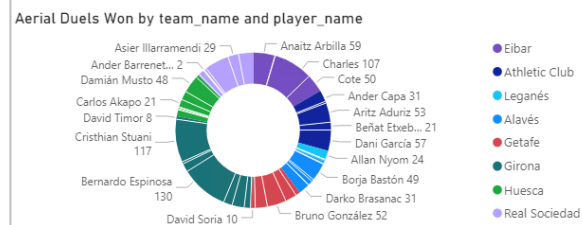
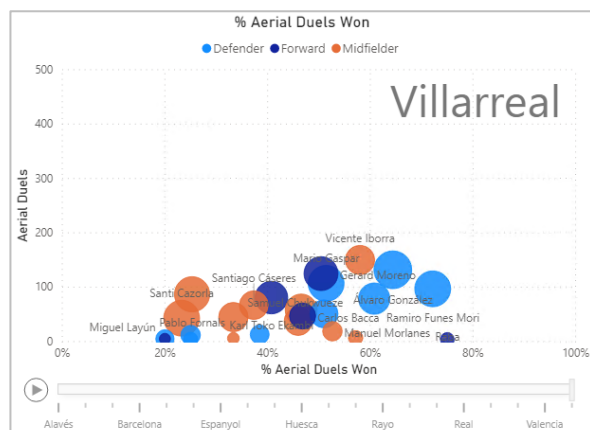


_ Mapas de Posesiones

_ Análisis de los duelos aéreos

_ Gráficos de radar con estadísticas de los máximos goleadores y asistentes de La Liga

_ Análisis de los últimos 5 clasificados

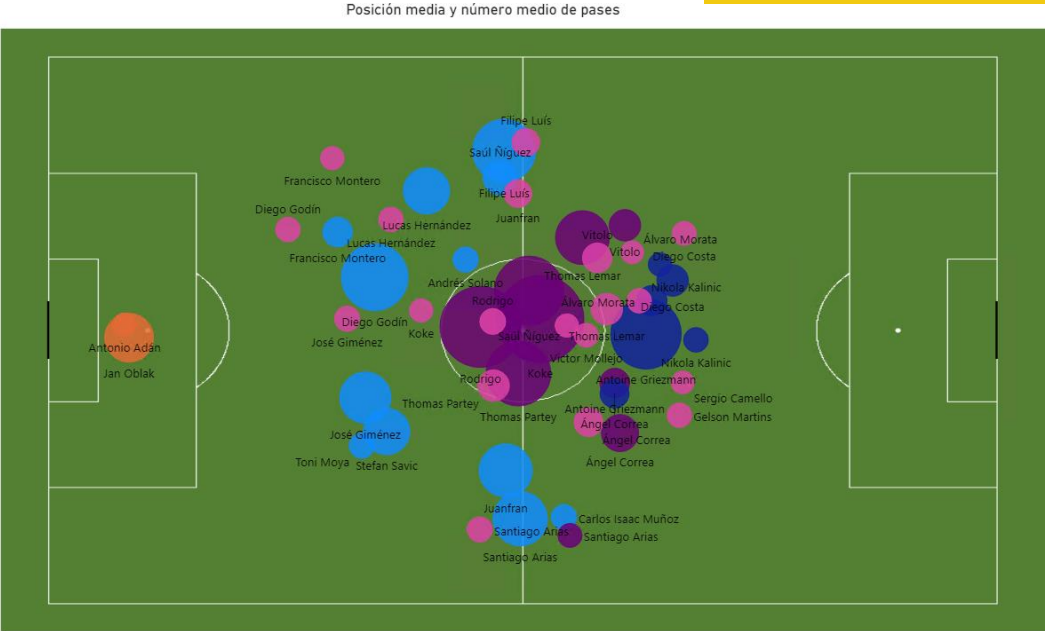
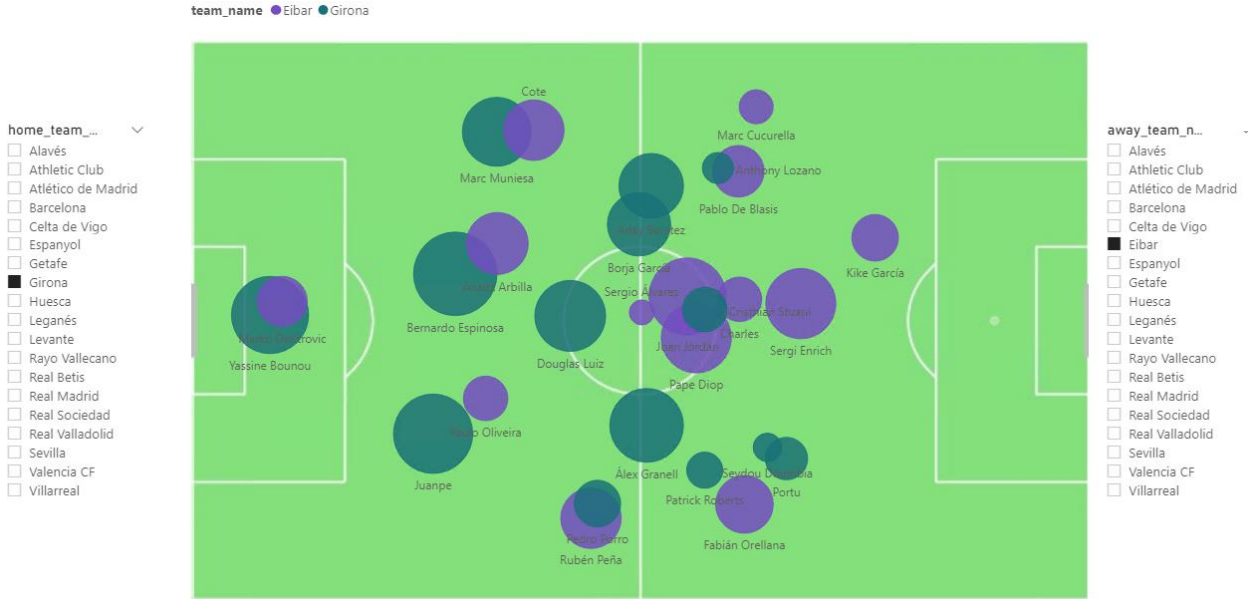


VISUALIZACIÓN DE RESULTADOS

POWER BI (*Pass Matrix.pbix*)



- _ Matriz de pases
- _ Posición media y número de pases



Pass Matrix

player_name	Arthur	Arturo Vidal	Coutinho	Gerard Piqué	Ivan Rakitic	Jordi Alba	Lionel Messi	Nélson Semedo	Sergi Roberto	Serg
Ivan Rakitic	100	83	115	223	2	255	241	139	277	
Sergio Busquets	169	153	133	217	253	165	283	92	138	
Gerard Piqué	115	115	39	2	274	62	86	210	265	
Jordi Alba	211	81	353	41	199		133	16	27	
Lionel Messi	85	116	147	42	167	135	2	67	127	
Sergi Roberto	42	132	37	161	248	18	218	54		
Arturo Vidal	49	1	31	139	78	77	178	127	80	
Clément Lenglet	63	34	68	137	141	217	11	24	16	
Coutinho	87	43		27	106	251	122	22	43	
Nélson Semedo	19	151	11	174	149	8	91		45	
Marc-André ter Stegen	41	35	21	208	60	68	23	37	69	
Luis Suárez	31	38	61	15	70	87	135	34	62	
Samuel Umtiti	23	39	67	97	81	112	10	12	37	
Ousmane Dembélé	23	45	50	15	70	93	99	42	26	
Malcom	11	29	5	3	8	13	14		24	

CONCLUSIONES Y LÍNEAS FUTURAS

CONCLUSIONES

_Optimización del modelo de datos

- Data mart
- Disponibilidad de cargar datos de otras temporadas y/o competiciones en un 'único repositorio'

_Procesos de carga

- Reducción de tareas manuales
- Pentaho Data Integration (PDI)
 - Fácil e intuitivo
 - Posibilidad de introducir mejoras

_Visualización

- Mejora notable al explotar la información estructurada procedente de la base de datos en vez de cargar ficheros XML directamente en Power BI

LÍNEAS FUTURAS

_Opción 1

- Extender y automatizar la interpretación de eventos del fichero F24
- Incorporar más OPTA feeds
 - F1-F75, especial hincapié en tracking (F53-55)

_Opción 2

- Integrar fuentes de datos procedentes de otros proveedores deportivos



