

# Construcció i explotació d'un magatzem de dades per a l'anàlisi d'informació sobre allotjaments turístics

## **Memoria**

Ferrán Casasús Rodó  
Enginyeria Tècnica d'Informàtica de Sistemes  
Consultor: Carles Llorach Rius  
Treball Fi de Carrera  
Curs 2012-2013 segon semestre

## Resumen

Trabajo centrado en el estudio y aplicación práctica de las técnicas conocidas bajo el nombre de “Almacenes de Datos” y concretamente en sistemas **ROLAP** para dar respuesta a las necesidades de los entornos decisorios. En él se detalla los requisitos para desplegar un “Análisis de información sobre establecimientos turísticos”, se expone la planificación, método basado en UML, el análisis, diseño, los pasos de recuperación de los datos, para su posterior transformación y carga en un base de datos relacional **MySQL** y la presentación de las consultas **OLAP** e informes.

## Índice de contenidos

1. Capítulo. Introducción.....	6
1.1. Justificación del proyecto.....	6
1.2. Objetivos del TFC.....	6
1.3. Enfoque y método .....	6
1.3.1. Arquitectura del proyecto.....	6
1.3.2. Infraestructura Informática.....	7
1.3.3. Método de trabajo y herramientas.....	7
1.4. Planificación.....	8
1.5. Tareas e hitos.....	8
1.5.1. Fase 1 Plan de trabajo.....	9
1.5.2. Fase 2 Análisis y diseño.....	10
1.5.3. Fase 3 Implementación.....	11
1.5.4. Fase 4. Presentación.....	11
1.6. Riesgos e incidencias.....	12
1.7. Productos Obtenidos.....	12
1.8. Descripción de los próximos capítulos.....	13
2. Capítulo. Análisis.....	13
2.1. Definición del contexto.....	13
2.1.1. Consultas requeridas.....	14
2.1.2. Fuentes de datos Disponibles.....	14
2.2. Criterios de un Almacén de Datos Decisional.....	15
2.3. Casos de Uso de la Explotación.....	16
2.3.1. Casos de uso Usuario.....	16
2.3.2. Caso de uso Constructor.....	17
2.3.3. Casos de uso Administrador.....	17
2.4. Modelo Conceptual Multidimensional .....	18

2.4.1. <i>Introducción al Modelo Conceptual</i> .....	18
2.4.2. <i>Criterios generales en el modelo conceptual</i> .....	18
2.4.3. <i>Nivel 1: modelo conceptual</i> .....	18
2.4.4. <i>Nivel 2: Paquete dimensión Comuns</i> .....	19
2.4.5. <i>Nivel 2: Paquete estrella Allotjaments</i> .....	20
2.4.6. <i>Nivel 2: Paquete estrella Demogrífics</i> .....	21
2.4.7. <i>Nivel 2: Paquete estrella Equipaments</i> .....	21
3. <i>Capítulo. Diseño</i> .....	22
3.1. <i>Modelo lógico multidimensional</i> .....	22
3.1.1. <i>Criterios generales para el Modelo Lógico Multidimensional</i> .....	23
3.1.2. <i>Notas sobre la representación del modelo lógico</i> .....	23
3.1.3. <i>Nomenclatura de los elementos</i> .....	23
3.1.4. <i>Modelo lógico nivel 1: Visión general</i> .....	23
3.1.5. <i>Nivel 2: Paquete lógico Comúns</i> . ....	24
3.1.6. <i>Nivel 2: Paquete lógico Allotjaments</i> . ....	24
3.1.7. <i>Nivel 2: Paquete lógico Demografic</i> .....	26
3.1.8. <i>Nivel 2: Paquete lógico Equipaments</i> . ....	26
3.2. <i>Modelo lógico del proceso ETL</i> .....	27
3.2.1. <i>Suposiciones generales aplicadas</i> .....	29
3.2.2. <i>Paquete ETL Comuns</i> .....	29
3.2.3. <i>Paquete ETL Allotjaments</i> .....	29
3.2.4. <i>Paquete ETL Demografics</i> .....	30
3.2.5. <i>Paquete ETL Equipaments</i> . ....	30
3.3. <i>Modelo Físico</i> .....	31
3.3.1. <i>Sistema Gestor de Base de Datos</i> .....	31
3.3.2. <i>Equipamiento</i> .....	32
4. <i>Capítulo Implementación</i> .....	34
4.1. <i>Exposición del trabajo de integración de los datos</i> .....	34
4.1.1. <i>Datos comunes</i> .....	36
4.1.2. <i>Demográficos</i> .....	37
4.1.3. <i>Alojamientos</i> .....	38
4.1.4. <i>Equipamientos</i> .....	41
4.2. <i>Proceso de creación de las consultas OLAP y de los informes</i> .....	42
4.2.1. <i>Proceso de creación de los metadatos y esquemas</i> .....	42
4.2.2. <i>Proceso de creación de los informes</i> .....	44
4.3. <i>Configuraciones y accesos</i> .....	44

5. Presentación de los informes realizados.....	44
5.1. Consultas OLAP.....	44
5.1.1. Consulta Establiments Comarques Ratios.....	45
5.1.2. Consulta Establiments Genere Turisme.....	47
5.1.3. Consulta Equipaments Comarca Ratios.....	49
5.1.4. Consultas OLAP de los data-marts.....	49
5.2. Informes.....	52
5.2.1. Informe Establiments Comarca (InformeEstablimentComarca.prpt).....	52
5.2.2. Informe Establiments Comarca Homes Dones Equipaments (InformeEstablimentsComarcaB.prpt) .....	54
5.2.3. Informe Equipaments Comarca (informeEquipamentsComarca.prpt).....	55
6. Conclusiones.....	57
7. Lineas de evolución futura.....	57

*Anexos:*

<i>Anexo: Por que ArgoUML?.....</i>	<i>59</i>
<i>Anexo Modelo Conceptual basado en UML.....</i>	<i>59</i>
<i>Anexo Modelo Lógico basado en UML.....</i>	<i>69</i>
<i>Anexo Modelo Lógico ETL basado en UML.....</i>	<i>76</i>

## **Índice de ilustraciones**

<i>Ilustración 1: Gantt: Plan de trabajo.....</i>	<i>9</i>
<i>Ilustración 2: Gantt Fase de análisis y diseño.....</i>	<i>10</i>
<i>Ilustración 3: Gantt Fase Implementación.....</i>	<i>11</i>
<i>Ilustración 4: Gantt Fase presentación.....</i>	<i>11</i>
<i>Ilustración 5: Caso de uso del Usuario.....</i>	<i>16</i>
<i>Ilustración 6: Caso de uso Constructor.....</i>	<i>17</i>
<i>Ilustración 7: Caso de uso Administrador.....</i>	<i>17</i>
<i>Ilustración 8: Modelo Conceptual.....</i>	<i>19</i>
<i>Ilustración 9: Paquete estrella Allotjament.....</i>	<i>20</i>
<i>Ilustración 10: Paquete estrella Demogràfic.....</i>	<i>21</i>
<i>Ilustración 11: Paquete estrella Equipaments.....</i>	<i>22</i>
<i>Ilustración 12: Nivel 1 Modelo lógico.....</i>	<i>24</i>
<i>Ilustración 13: Paquete lógico Allotjaments.....</i>	<i>25</i>
<i>Ilustración 14: Paquete lógico Demogràfic.....</i>	<i>26</i>
<i>Ilustración 15: Paquete Lógico Equipaments.....</i>	<i>27</i>



<i>Ilustración 16: Arquitectura Pentaho.....</i>	<i>33</i>
<i>Ilustración 17: Imagen del trabajo que funciona como guía de los pasos a realizar.....</i>	<i>35</i>
<i>Ilustración 18: Generación de la dimensión temporal.....</i>	<i>36</i>
<i>Ilustración 19: Generación de nombres de municipios.....</i>	<i>37</i>
<i>Ilustración 20: Transformación de datos demográficos.....</i>	<i>38</i>
<i>Ilustración 21: Transformación alojamientos.....</i>	<i>39</i>
<i>Ilustración 22: Transformación alojamientos descarte de filas.....</i>	<i>40</i>
<i>Ilustración 23: Transformación de equipamientos.....</i>	<i>41</i>
<i>Ilustración 24: La generación de los categorías de equipamientos.....</i>	<i>42</i>
<i>Ilustración 25: Ratios Establecimientos comarca lado izquierdo.....</i>	<i>46</i>
<i>Ilustración 26: Ratios Establecimientos comarca lado derecho.....</i>	<i>46</i>
<i>Ilustración 27: Volcado de la pantalla de la consulta establecimientos género lado izquierdo .....</i>	<i>47</i>
<i>Ilustración 28: Volcado de la pantalla de la consulta establecimientos género lado derecho</i>	<i>48</i>
<i>Ilustración 29: Volcado de pantalla con la consulta equipamientos comarca.....</i>	<i>49</i>
<i>Ilustración 30: El cubo de equipamientos.....</i>	<i>50</i>
<i>Ilustración 31: El cubo de Demográfico .....</i>	<i>51</i>
<i>Ilustración 32: El cubo de los establecimientos.....</i>	<i>51</i>
<i>Ilustración 33: Los totales a nivel año tipo y provincia establecimientos comarca.....</i>	<i>53</i>
<i>Ilustración 34: El detalle del informe establecimientos comarca.....</i>	<i>53</i>
<i>Ilustración 35: El detalle del informe Establiments Comarca Homes Dones Equipaments.</i>	<i>54</i>
<i>Ilustración 36: Totales provincia, categoría y año informe Establiments Comarca Homes Dones Equipaments.....</i>	<i>55</i>
<i>Ilustración 37: Cabecera y cuerpo del informe Equipaments Comarca .....</i>	<i>56</i>
<i>Ilustración 38: Los totales a nivel tipo y año del informe Equipaments Comarca .....</i>	<i>56</i>

## **Índice de tablas**

<i>Tabla 1: Plan de contingencias.....</i>	<i>12</i>
<i>Tabla 2: Tabla de tipos y categorías de Establecimiento.....</i>	<i>20</i>

## 1. Capítulo. Introducción

---

Para el trabajo de final de la carrera d'Enginyeria Tècnica d'Informàtica de Sistemes hay que desarrollar un trabajo de síntesis aplicado en un área en concreto, en este caso es el análisis de las técnicas actuales para diseñar bases de datos para un almacén de datos (Data Warehouse), su consulta y manipulación.

### 1.1. JUSTIFICACIÓN DEL PROYECTO.

El proyecto de “Análisis de información sobre establecimientos turísticos” esta enfocado a plantear un reto con la intención de ofrecer un ejemplo de los problemas más comunes en el desarrollo de “Almacenes de datos”.

### 1.2. OBJETIVOS DEL TFC

En una sola frase y tal como se expone en el enunciado del trabajo el objetivo principal será poseer experiencia en el diseño, construcción y explotación de un almacén de datos a partir de la información disponible en una base de datos transaccional.

Si lo detallamos un poco más, tenemos:

- Conocer e investigar en las técnicas actuales para diseñar bases de datos para un almacén de datos según los principios de los almacenes de datos físicos **ROLAP**. Estar en situación para considerar factores tales como desnormalización de tablas, inclusión de información agregada, históricos, etc.
- Comprender los problemas derivados de la recuperación de datos de múltiples fuentes y formatos distintos y duplicados.
- Aplicar un criterio metodológico para todo el proceso y utilizando una metodología integradora como el **UML** adaptado a los almacenes de datos.
- Estudiar un conjunto de herramientas concretas que son necesarias para crear y gestionar un sistema de información decisorio basado en almacenes de datos.
  - Sistema ETLs: **Spoon/Keettle**.
  - Servidor de Inteligencia Empresarial (Business Intelligence). **Pentaho BI Server**
  - Sistemas de consultas: **Pentaho Report Designer** y análisis OLAP con **jPivot** y **SAIKU** sobre el servidor **Mondrian**

### 1.3. ENFOQUE Y MÉTODO

#### 1.3.1. ARQUITECTURA DEL PROYECTO

Este es un proyecto para los sistemas de información decisorios. Se entiende como sistema de información decisorios, a aquella parte del SI global de la organización dedicada a dar soporte a los diferentes tipos de procesos de toma de decisiones.

Se basará en software abierto de contrastada eficacia que cubrirá las 3 grandes capas: la importación de datos, su almacenamiento y posterior consulta. A un nivel de mas detalle tenemos:

- Almacén de datos, que está orientado a temas o conceptos , con información de largos periodos de tiempo, no volátil e integrado.
- Metadatos, que describen cual es la estructura de los datos que se almacenan y como se relacionan.
- Sistemas ETL "Extract, Transform, and Load" que lee la información de los datos primarios (fuentes de datos operativas), realiza procesos de transformación para ser almacenados (filtrado, adaptación, cambios de formato) y finalmente guardarlos.
- Sistema de consultas y explotación de usuario. Que permite de una forma ágil suministrar las consultas que precisa el usuario.

El detalle de la pila de herramientas a utilizar se estudiará en la fase de diseño.

### 1.3.2. INFRAESTRUCTURA INFORMÁTICA.

La infraestructura informática se compone de 2 grandes áreas, por un lado están las herramientas para la documentación y presentación del proyecto y por otro las tecnologías necesarias para la implantación y explotación del producto a elaborar.

- Para el caso del Pla de Treball se ha utilizado un equipo en windows 97 con OpenOffice.org 3.3.0 para los documentos. Notepad++ v6.2.2 para notas y borradores. Ganttproject 2.0.10 para realizar los diagramas de Gantt y Argo UML para el como herramienta de diagramación del Análisis y Diseño (ver anexo sobre ArgoUML).
- Para caso del producto a elaborar, se realizará en la Maquina Virtual suministrada por la UOC bajo VirutalBox. Con ella se simulara el entorno de desarrollo y de usuario final.

### 1.3.3. MÉTODO DE TRABAJO Y HERRAMIENTAS.

Para presentar y desarrollar la fase de análisis se ha trabajado bajo el modelo *Unified Modeling Language (UML)* con la herramienta *Argo UML* [ref-argo],

#### Por que UML?

El utilizar el **UML** como sistema de diagramación y repositorio implica que había que tomar una decisión sobre como realizar el modelado: si se debía de utilizar el modelo *Entidad/Relación (ER)* y adaptarlo a la representación **UML** o utilizar una propuesta dentro del universo **UML**. Me he inclinado por lo segundo a partir de las siguientes lecturas:

[ref-prat-akoka] Que indica como transponer la visión de los 3 modelos utilizados en las bases de datos (conceptual, lógico y físico) a los almacenes de datos. Los autores argumentan que mucha de la literatura no aparece clara dicha distinción y ellos creen que es útil y necesaria que exista.

[ref-LujanS] Tesis de Sergio Luján que presenta y detalla una extensión del *Unified Modeling Language (UML)*, para desarrollar según el método *Unified Process (UP)* y aprovechando el uso de estándares en el desarrollo de los almacenes de datos. Luján S. expone un sistema de diagramación que cubre todas las partes del desarrollo de un almacén de datos.

[ref-LujanS] defiende la decisión de utilizar el UML por cinco consideraciones básicas:

- *“UML sigue el paradigma OO, y que se ha demostrado tener una semántica mas rica y que otros paradigmas porque OO modelos (pueden) estar mas cerca a la concepción del usuario”*
- *“EL UML es un lenguaje conocido por los ingenieros informáticos”*
- *“El uso de los perfiles facilita el diseño, y así los desarrolladores de Data Warehouse DW pueden utilizar los conceptos que ahora ya aplican. Esto presenta la ventaja que se puede utilizar el UML para DW sin ser un experto en UML”*
- *“El UML es el estándar de Object Management Group (OMG) y unifica muchos años de esfuerzo en el análisis y diseño OO”*
- *“UML se ha convertido en el modelo dominante en el mundo OO”*

## Como se va presentar el UML en la memoria?

El UML es un lenguaje **visual** utilizado como método sistemático para definir los aspectos de del análisis y diseño de un producto, pero para la memoria se precisa más bien un **relato** de los realizado y los gráficos son de soporte y comprensión del texto.. Para solucionar esta dicotomía entre los dos formatos de exposición, se ha decidido, dejar toda la documentación generada con el UML como anexos a la memoria. Así tenemos que en el anexo (Anexo Modelo Conceptual basado en UML) hay una exposición en detalle de los artefactos utilizados y los criterios de uso en el modelo UML.

### **1.4. PLANIFICACIÓN**

Para conseguir una correcta implementación del almacén de datos corporativo para ONdO, será construido por etapas, en cada uno de ellas se irá solucionando una problemática concreta.

Cada final de etapa se verá validada por la entrega de la PAC con el material elaborado en dicha fase.

Así tenemos 4 grandes fases:

- Plan de trabajo y análisis preliminar de requerimientos.
- Análisis de requerimientos y diseño conceptual y técnico
- Implementación
- Presentación

### **1.5. TAREAS E HITOS**

Dentro de las tareas e hitos he tenido en cuenta, la necesidad de formación en tres ámbitos: en los almacenes de datos, en las tecnologías que los soportan, en la elaboración del proyecto.

Respecto a la asignación de jornadas y debido a mi trabajo como freelance, y por lo tanto dependiente de los plazos de entrega de mis trabajos con los clientes, puede implicar ciertos desajustes en la planificación, en todo caso, intentaré dar prioridad al TFC. La planificación está pensada en bloques semanales con jornadas de 4 horas.

La transcripción del diagrama de Gantt al Plan de Trabajo lo he realizado por secciones para que se pueda apreciar correctamente los días y jornadas asignadas a cada tarea.

Leyendas:

- Los hitos aparecen como rombos. En rojo los hitos de entrega.
- Los plazos de color ■ indican tareas documentación.

### 1.5.1. FASE 1 PLAN DE TRABAJO.

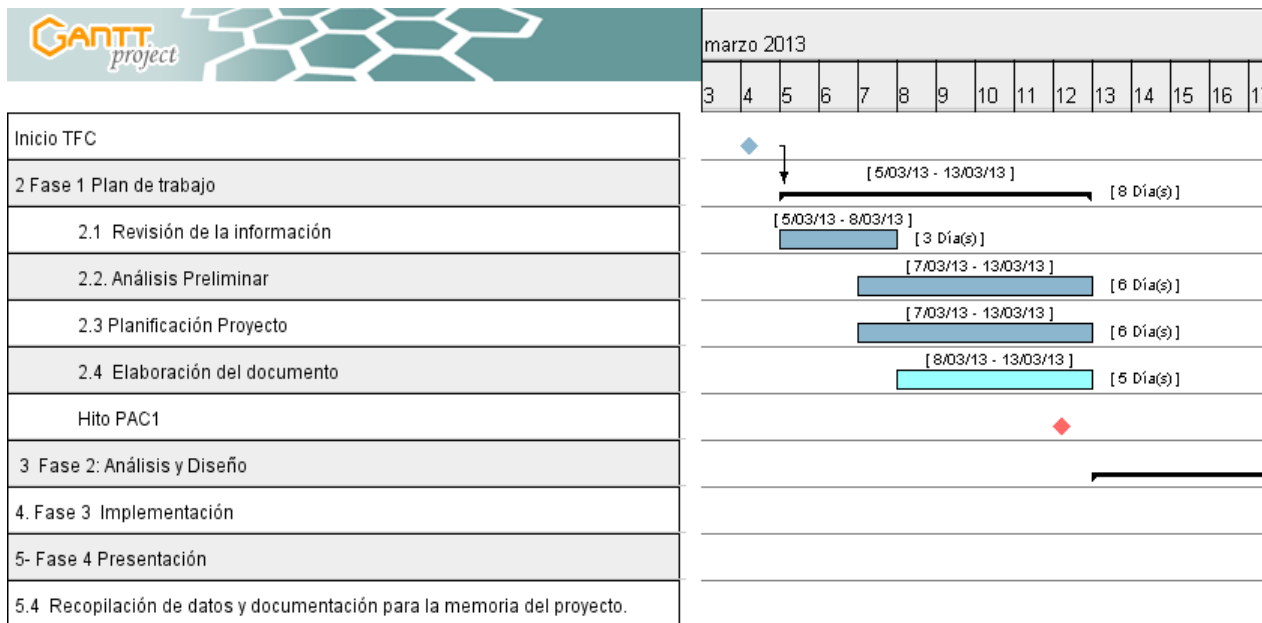


Ilustración 1: Gantt: Plan de trabajo

Como era de esperar esta fase se ha cumplido sin mayores desviaciones.

### 1.5.2. FASE 2 ANÁLISIS Y DISEÑO.

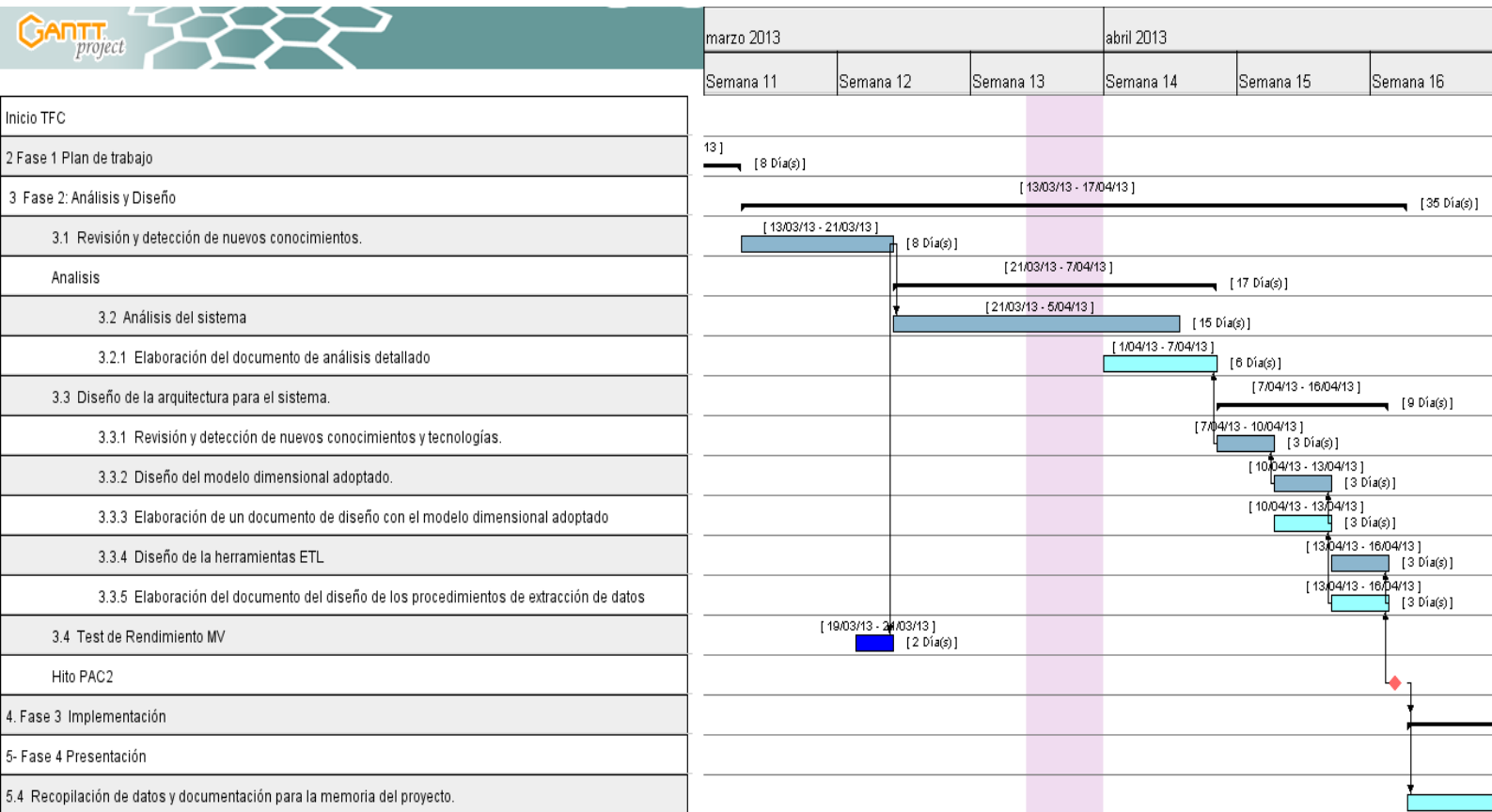


Ilustración 2: Gantt Fase de análisis y diseño

Al ser la fase con un desglose de tareas mayor, el proceso de cálculo de los plazos, se ha realizado empezando desde la fecha de entrega para ir asignando jornadas a las tareas en el orden inverso al que serán ejecutadas.

El plan se ha visto alterado ya que por razones de carga de trabajo, se tuvo que concentrar todo el proyecto en las tres últimas semanas, también se desvió el volumen de horas dedicado a formación ya que también implicó investigar en cual sería el mejor uso de la diagramación UML aplicado en toda la fase, por suerte el tiempo previsto para el análisis fue menor, y en cierta medida compenso. Como el análisis y diseño del proyecto se trabajo previamente con diagramas UML, incluyendo el diseño ETL, la redacción del documento quedó relegada al final de la fase.

### 1.5.3. FASE 3 IMPLEMENTACIÓN

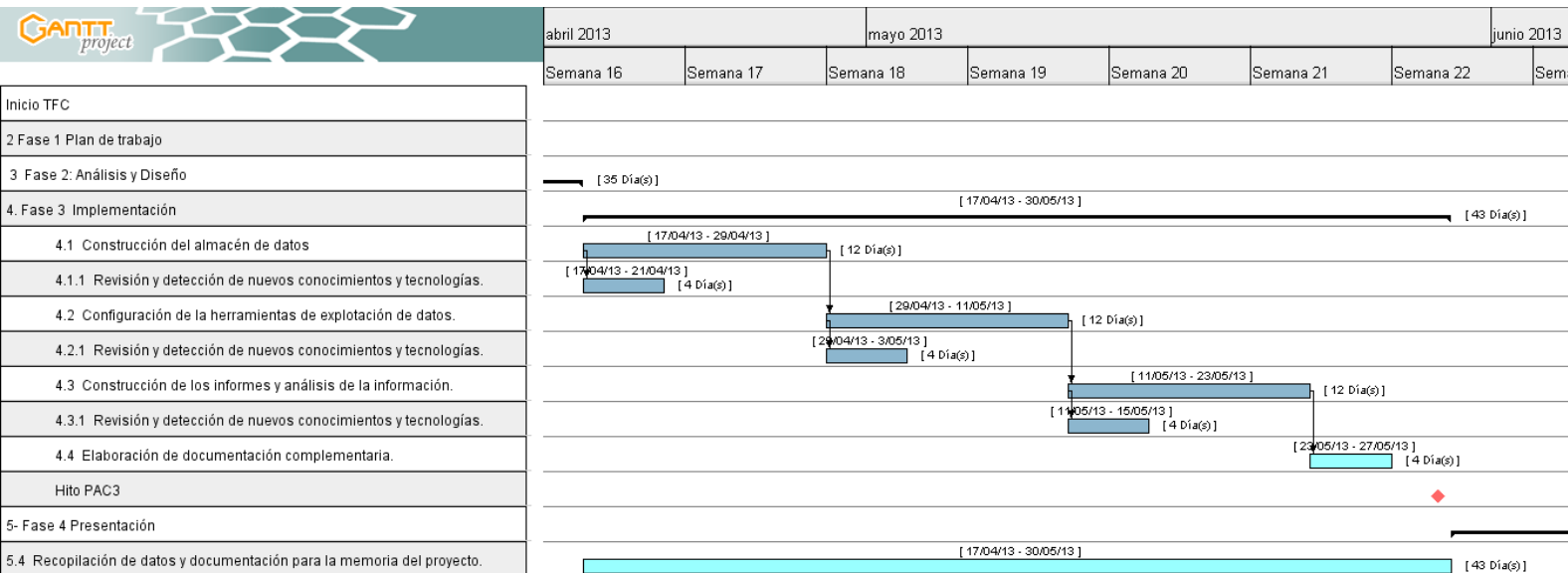


Ilustración 3: Gantt Fase Implementación

En esta fase, es donde ha habido, la mayor desproporción de horas y plazos, por otra parte normal ya que no se conocían todas las tecnologías que se utilizan, por lo que se realizaron varias aproximaciones, antes de encontrar el uso correcto de las mismas.

### 1.5.4. FASE 4. PRESENTACIÓN

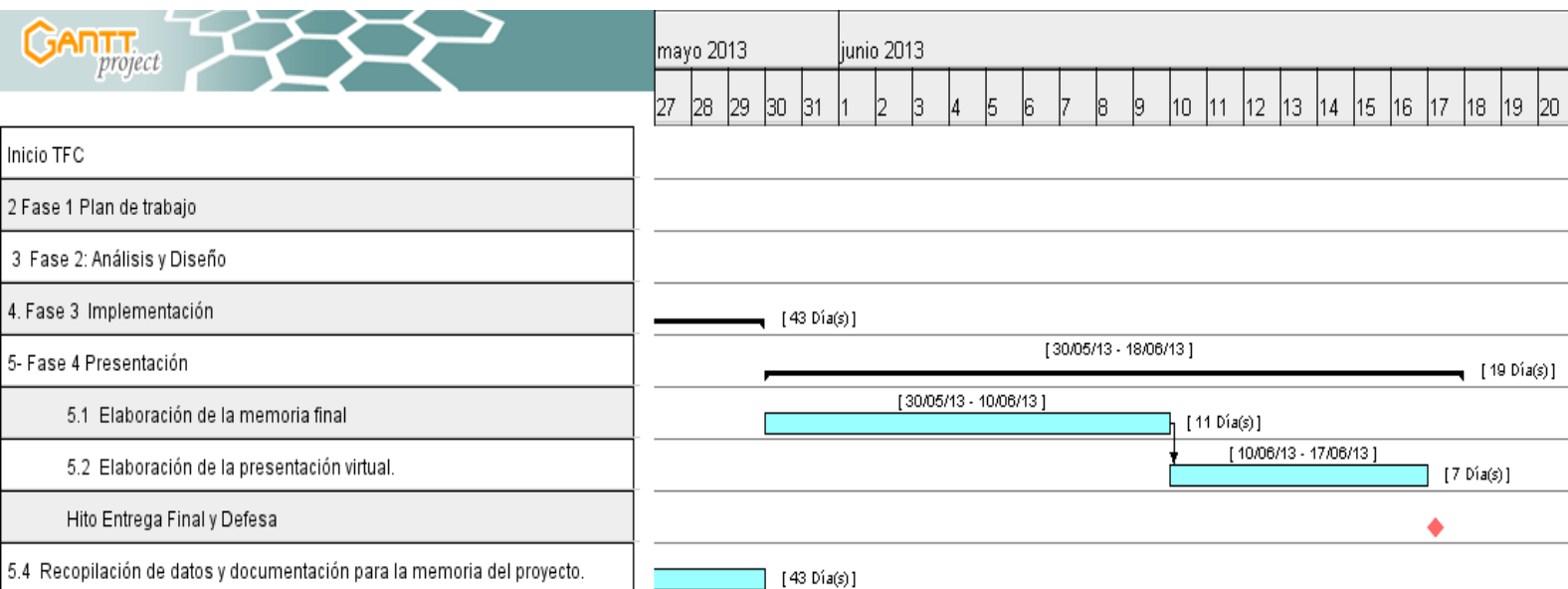


Ilustración 4: Gantt Fase presentación

## 1.6. RIESGOS E INCIDENCIAS

En la siguiente tabla se relacionan las incidencias mas destacadas que se pueden dar durante el desarrollo del Proyecto de Final de Carrera, con sus planes de contingencia.

Tabla 1: Plan de contingencias

• Incidencia	• Plan de Contingencia.
<ul style="list-style-type: none"> <li>• El equipo no aguanta la maquina virtual</li> </ul>	<ul style="list-style-type: none"> <li>• Posibilidad de ampliar memoria.</li> <li>• Posibilidad de crear disco virtual exclusivo para memoria intermedia.</li> </ul>
<ul style="list-style-type: none"> <li>• Fallos en la maquina virtual</li> </ul>	<ul style="list-style-type: none"> <li>• Recuperar de las copias a disco duro externo de la Maquina virtual</li> </ul>
<ul style="list-style-type: none"> <li>• Perdida de ficheros de documentación del proyecto fuera de la maquina virtual</li> </ul>	<ul style="list-style-type: none"> <li>• Recuperar de las copias a disco duro externo.</li> </ul>
<ul style="list-style-type: none"> <li>• El desarrollo lleve a un camino sin salida, y es necesario deshacer lo hecho.</li> </ul>	<ul style="list-style-type: none"> <li>• Recuperar de versiones anteriores.</li> </ul>
<ul style="list-style-type: none"> <li>• Problemas de salud</li> </ul>	<ul style="list-style-type: none"> <li>• Si el tema es grave, no hay otra salida que parar el proyecto y dejarlo para el próximo semestre.</li> <li>• Solo hay margen de recuperación si después se amplia la disponibilidad horaria.</li> </ul>
<ul style="list-style-type: none"> <li>• Reducción de la disponibilidad horaria.</li> </ul>	<ul style="list-style-type: none"> <li>• Si la reducción es aguda no hay otra salida que parar el proyecto y dejarlo para el próximo semestre.</li> <li>• Solo hay margen de recuperación si después se amplia la disponibilidad horaria.</li> <li>• En todo caso el TFC tendrá prioridad sobre posibles trabajos que representen una disminución de la disponibilidad horaria.</li> </ul>

## 1.7. PRODUCTOS OBTENIDOS

- El fichero **UML** expuesto de forma expositiva en los anexos.
- El fichero **Spoon** para llevar a cabo el ETL que está en la máquina virtual
- Los informes
  - EstablimentsComarca
  - Establiments Comarca Homes Dones
  - Equipaments Comarca



- Las consultas OLAP
  - Ratios entre Establecimientos y Comarcas
  - Ratios entre Establecimientos habitantes por género y por grupo de equipamientos relacionados directamente con el turismo.
  - Ratios entre Equipamientos y Comarcas.
  - Consultas preparadas para generar nuevas consultas sobre los datos.
- Una maquina virtual en Windows que contiene tanto las herramientas cliente, como servidor
- Base de datos MySQL para el *Stagging Area* los procesos de transformación. Alojado en la máquina virtual
- Base de datos MySQL para el Proyecto Alojado en la máquina virtual

## **1.8. DESCRIPCIÓN DE LOS PRÓXIMOS CAPÍTULOS.**

En los próximos capítulos se desarrolla el análisis del proyecto (el modelo conceptual) presentando los grandes grupos de datos que de forma transversal se repiten en todas las fase de desarrollo. Se sigue con el diseño (el modelo lógico), hay un apartado dedicado al proceso ETL desde un punto de vista lógico (es decir enfocado a analizar la estructuras de las fuentes de datos y sus transformaciones para ser cargados en el almacén de datos), a continuación se relaciona de tablas creadas (el modelo físico), se planeta como se hizo la instalación y finalmente se presentan la lista de aplicaciones y informes generados.

## **2. Capítulo. Análisis**

---

El proceso de análisis se ha centrado en los requerimientos actuales, pero a su vez se ha intentado dejar abierto el sistema a futuras necesidades, sin que estas no significasen una complicación mayor en el sistema<sup>1</sup>.

### **2.1. DEFINICIÓN DEL CONTEXTO**

Se requiere desarrollar un sistema de almacén de datos, para un estudio comparativo y evolutivo de los diversos tipos de establecimientos de alojamiento turístico, comparándolo con la zona geográfica, su dimensión, su población y con la oferta de equipamientos de la zona. Esto permitirá tener indicadores sobre el impacto de la industria de alojamiento turístico en una zona determinada.

En un principio, con la lectura de los requerimientos, puede entenderse que va ser un producto a ser utilizado de manera puntual en un proyecto de consultoría más amplio.

En todo caso, y analizando la información de las fuentes de datos, creo que si en vez de realizar un análisis puntual, el verdadero valor de las correlaciones que se estudian en este proyecto, estaría en un seguimiento durante los próximos años , y por lo tanto ir incorporando tanto información histórica, como nueva en el almacén de datos.

Y es sobre la anterior visión como se ha analizado y diseñado tanto el almacén de datos, como

---

<sup>1</sup>Dicha restricción parte de la idea de ver el proyecto, no solo como trabajo de final de carrera, sino también como el proyecto para unos usuarios reales, por lo tanto han de considerar el coste que supone tener un sistema mas o menos flexible a modificaciones futuras.

los procesos **ETL**

### 2.1.1. CONSULTAS REQUERIDAS.

El número de informes que se requieren, son como mínimo los definidos en el enunciado, un total de 10, lo cuales los he agrupado en dos grandes grupos en función de las características de navegación

Hay un primer grupo que comparan establecimientos con zona geográficas y su población (habitantes y dimensión) , en todos ellos ha de poderse navegar por año y por tipo de establecimiento y agrupación geográfica.

*”Total d’establiments”.*

*”Total de places”.*

*”% de places respecte població”.*

*”Oferta mitjana de places”.*

*”Indicador d’establiments vs habitants per gènere”.*

*”Indicador de places vs persones”.*

*”Quantitat de places ofertes / superfície del territori”.*

Tenemos un segundo grupo que compara las relaciones con los equipamientos, en esos casos ha de poderse, además, navegar por el árbol de tipos de equipamientos.

*”Nombre d’establiments/Nombre d’equipaments”*

*”% de població per equipament”*

*”Indicador d’equipaments vs població”.*

En todo caso, el sistema debería de permitir cualquier cruce geográfico y temporal entre la dimensión de establecimientos hoteleros y sus valores número y plazas, la dimensión de población y sus valores de número de habitantes y el área y la dimensión de equipamientos con sus valores de cantidad y tipos.

Lo anterior implica que hay que añadir la dimensión tiempo en equipamientos, lo que permitirá realizar importaciones anuales para que reflejen los cambios se produzcan en los equipamientos.

### 2.1.2. FUENTES DE DATOS DISPONIBLES

Las fuentes de datos que disponemos no son datos operacionales del sistema sino informes elaborados desde distintos orígenes, principalmente públicos.

Las fuentes proveen valores durante series de años, y se ha observado que en algunas series los nombres identificativos varían, así como también varían las categorías que incluyen.

Como criterio de diseño se ha establecido que siempre que sea posible utilizar la definición de nombres y categorización de los valores del años 2012.

Existen 3 diferentes fuentes principales:

- *poblacio.csv* con información sobre el número de habitantes entre los años 2012 y 2006 y su extensión en km2, además con el detalle en el último año 2012 del número de hombres y mujeres. Cada población tiene su código INE.

- *equipaments.csv* con información a fecha de 31/12/2012 sobre el nombre, dirección, municipio/población, comarca, código postal, teléfono, longitud, latitud, categorías (multi-nivel, por ejemplo “Equipaments | Cultura | Arxius”) y localización.
- *Establecimientos año .txt.*, Los años son del 2006 al 2012. Contiene información proveniente de dos fuentes: *Federació Catalana d’Allotjaments Turístics* y *Idescat*, a partir de los datos del Departament d’Empresa i Ocupació
  - Fuente: **Federació Catalana d’Allotjaments Turístics**: número de establecimientos y sus plazas, en formato desglosado y total por
    - Distribución geográfica: *Catalunya, àmbits, províncies, comarques*
    - Tipos: *Hotels, Càmpings, Turisme Rural*
    - Hoteles desglosado por *Hotels, pensions o estrelles oro, estrella argent,*
    - Càmpings desglosado por categorías: *luxe, 1ª, 2ª, 3ª*
    - Turismo Rural desglosando por: *Casa de poble compartida, independent, Masia, Masoveria*. Existen variaciones en las tipificaciones entre años.
  - Fuente: **Idescat**, a partir del **Departament d’Empresa i Ocupació**. solo año 2012 que da información sobre plazas por
    - Distribución geográfica: *Catalunya, àmbits, províncies, comarques*
    - Tipo de Turismo Rural: *Casa de poble compartida, independent, Masia, Masoveria*.

## 2.2. CRITERIOS DE UN ALMACÉN DE DATOS DECISIONAL

Para poder dar solución a los requisitos del proyecto es necesario una herramienta que permita hacer consulta en línea y de una forma dinámica de la información (**OLAP**) lo que implica que el sistema de consulta y de almacenamiento ha de cumplir una serie de criterios

- Toda la información está centralizada y accesible bajo el sistema **OLAP**
  - La información se puede actualizar al instante, el sistema ha de ofrecer la posibilidad de tener automatizado el proceso de extracción, transformación de los datos (**ETL**).
  - El acceso a los datos ha de ser lo máximo de rápido. La estructura de la información se realiza en forma de estrella, se desnormalizan los datos, se diseñan sistemas cache para guardar información agregada.
  - El sistema de almacenamiento ha de poder soportar altos volúmenes de datos.
  - El acceso a los datos por parte del usuario ha de ser fácil de entender y manipular
- Existe otro tipo de requisitos que son dependientes del diseño del problema
- Todos los datos históricos han de estar disponibles, es decir que se pueda seguir la cronología de un dato con sus sucesivas transformaciones.
  - La definición de cada dato es única, no existen varios nombres para referirse a ella.
  - Los datos están estandarizados, solo existe una definición y un conjunto de valores por cada dato.

Respecto al almacenamiento se va a utilizar una herramienta **OLAP** sobre base de datos relacional (**ROLAP**). existen otros modelos como en base de datos especializadas para contener almacenes de datos que son conocidas con el nombre (**MOLAP**) e incluso existen casos híbridos que combinan ambos mundos (**HOLAP**)

### 2.3. CASOS DE USO DE LA EXPLOTACIÓN.

Se han identificado los siguientes tipos de usuarios genéricos: el que consulta y navega con las herramientas de consultoría, el constructor de informes adicionales y finalmente quien administra el sistema y realiza la importación de datos

- Usuario: Utiliza el portal para consultar los informes, abre los listados predefinidos y aplica filtros. Sería el caso de quién accede a los informes definidos en los requerimientos, para después navegar por las diversas dimensiones. Además de los informes podrá manipular los cubos definidos en el sistema.
- Constructor: Puede crear informes adicionales y nuevas vistas tanto con las herramientas simples de usuario que dispone el sistema **Web Based Ad Hoc Query and Reporting**. como con herramientas mas sofisticados utilizando el **Report Designer** y crear nuevos modelos OLAP **Modrian (Pentaho)**.
- Administrador: Responsable del servidor, y de validar el proceso de ETL, con lo que no solo tiene que tener un perfil técnico, si no, conocer las fuentes de datos, su estructura y estar preparado para detectar fuentes erróneas.

#### 2.3.1. CASOS DE USO USUARIO

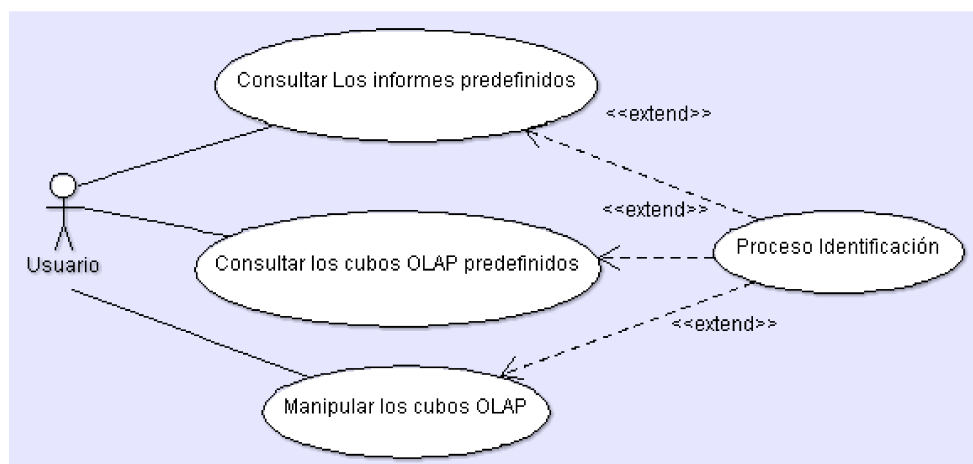


Ilustración 5: Caso de uso del Usuario

(Usuario) Consultar los informes predefinidos.

El usuario tras identificarse elegirá uno de los informes predefinidos, entrando los filtros que el informe haya definido y el sistema devolverá el informe en el formato de documento que el usuario haya elegido entre una lista.

(Usuario) Consultar los cubos OLAP predefinidos.

El usuario tras identificarse elegirá uno de los cubos que el sistema le ofrecerá a escoger, tras la

selecciona, el sistema le mostrar el cubo, pudiendo el usuario navegar por las diferentes jerarquías y miembros.

(Usuario) Manipular los cubos.

Se ha incluido ese caso por la existencia de una herramienta interactiva y “usable”, permite al usuario avanzado manipular unas fuentes de datos. El usuario seleccionara uno de los cubos disponibles y generara la consulta.

El usuario podrá guardar la consulta para reutilizar en posterior ocasiones.

Restricción esté caso solo estará permitido a los estrellas que estén definidos en el sistema, actualmente son demográfico/territorial, establecimientos y equipamientos

### 2.3.2. CASO DE USO CONSTRUCTOR

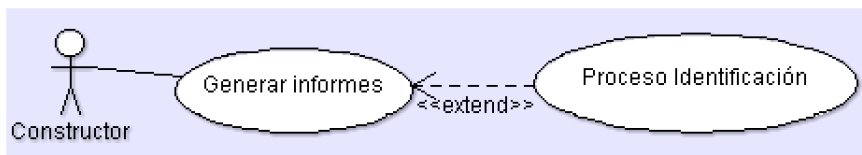


Ilustración 6: Caso de uso Constructor

(Constructor) Generar nuevos informes.

El constructor tras identificarse accederá al sistema para utilizar tanto un herramienta de generación de informes interactivos como una que permita documentos mas complejos..

### 2.3.3. CASOS DE USO ADMINISTRADOR

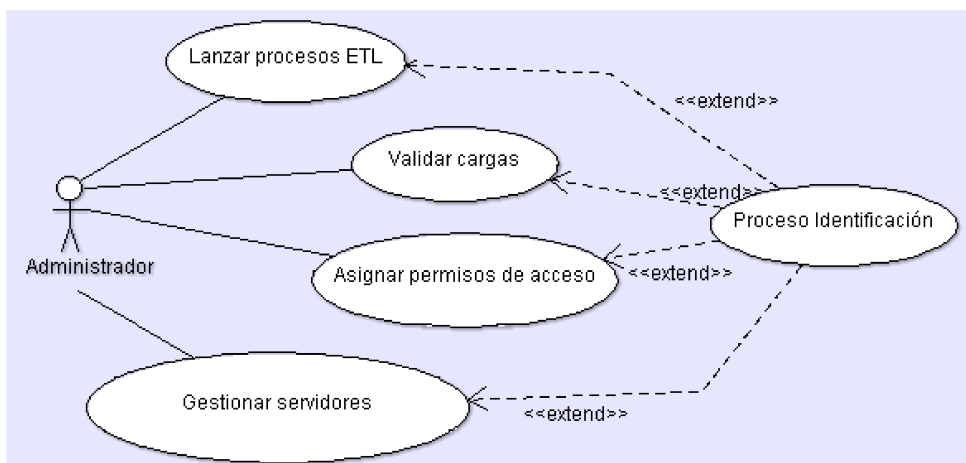


Ilustración 7: Caso de uso Administrador

(Administrador) Lanzar y controlar proceso ETL

El administrador tras identificarse accede al gestor ETL para realizar un proceso de transformación y carga, el sistema ha de mostrar por defecto la lista de tareas, de la cual existirá una que englobara a todas las otras. El administrado modificara los parámetros necesarios y lanzara

cada una de las tareas, el sistema ha de avisar de la entrada de incoherencia y del número de registros actualizados. En caso de incoherencia el sistema ha de parar el proceso.

#### (Administrador) Validar

El Administrador tras el proceso de ETL accede a la herramienta de análisis OLAP y de una manera interactiva podrá configurar un esquema OLAP para poder validar los resultados con los datos de entrada, mediante catas aleatorias de datos.

#### (Administrador) Asignar permisos de acceso a servidor BI

El administrador tras identificarse tendrá la opción de asignar a los documentos e informes creados los permisos de acceso para los usuarios finales.

#### (Administrador) Gestionar los servidores.

El administrador tras identificarse en una consola gestionara las programación de tareas, las configuraciones de los servidores y conexiones a las base de datos, el mantenimiento de los servicios.

## **2.4. MODELO CONCEPTUAL MULTIDIMENSIONAL**

### **2.4.1. INTRODUCCIÓN AL MODELO CONCEPTUAL**

En síntesis y tal como se expresa [ref-LujanS] podemos definir la necesidad del modelo conceptual

- Entender el ámbito o dominio real en el que esta insertado el problema
- Dar razón sobre dicho ámbito
- Conseguir sincronizar la visión del problema de todos los implicados en el proyecto.

### **2.4.2. CRITERIOS GENERALES EN EL MODELO CONCEPTUAL**

Guardar toda la información correcta que provea las fuentes de datos

Tener la información al mas bajo nivel de detalle que suministran las fuentes de datos

Nada se borra y se refleja su historia de cambios.

### **2.4.3. NIVEL 1: MODELO CONCEPTUAL**

En el apartado de “Definición del contexto” se han identificado dos grandes grupos de consultas que ha de responder el sistema:

- Comparativas de establecimientos con zona geográficas y su población (habitantes y km<sup>2</sup>) , en todos ellos ha de poderse navegar por año y por tipo de establecimiento y agrupación geográfica.
- El sistema a de soportar comparativas y ser capaz de establecer relaciones entre los establecimientos y los equipamientos, Además ha de poderse , navegar por el árbol de tipos

de equipamientos y clases de equipamientos.

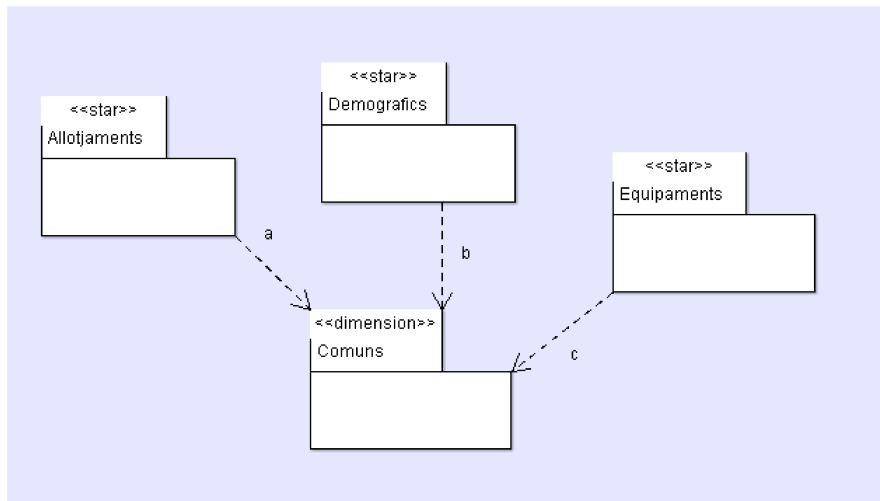


Ilustración 8: Modelo Conceptual

Se ha dividido el análisis en 4 grandes grupos o **paquetes UML** con lo que tenemos 3 formaciones de estrellas (paquetes <<start>>): *Allotjaments*, *Demografics*, *Equipaments*. Y una dimensión compartida por las tres estrellas llamada *Comuns*. Así tenemos en un primer nivel una visión de los diferentes tipos de valores que conforma el producto.

- El paquete *Allotjaments* contiene la estrella con los valores sobre el alojamiento hotelero y su dimensión de categorías de establecimientos y las dimensiones comunes: *Any* y *Comarca*
- El paquete *Demografics* contiene la estrella con los valores sobre la demografía y extensión geográfica, además de las dimensiones compartidas *Any* y *Municipi*
- El paquete *Equipaments* contiene la estrella con los valores agregados de equipamientos y su dimensiones de categorías de equipamientos (*Equipament*) e instalaciones (*Instalacions*), además de las dimensiones compartidas *Any* y *Municipi*

#### 2.4.4. NIVEL 2: PAQUETE DIMENSIÓN COMUNS

Dimensión temporal *Any*:

En este caso tenemos, a la dimensión *Any* que solo tiene el nivel base de *Any*. Este es el -único caso en donde la clave no será asignado automáticamente por la base de datos, sino será el año de cuatro cifras. A su vez será utilizado como descriptor en el sistema OLAP. Aparte del año, se ha previsto un elemento más con la intención de permitir filtrar por antigüedad, es decir por año

Dimensión territorial *Municipi*:

Tiene como base de más bajo nivel al municipio (*Municipi*) que después asciende por *Comarca* que a su vez se puede ascender por *Ambit* y *Provincia*.

Dimensión territorial *Comarca*

Tiene como base de más bajo nivel a la *Comarca* desde la que se puede ascender por *Ambit* o *Provincia*.

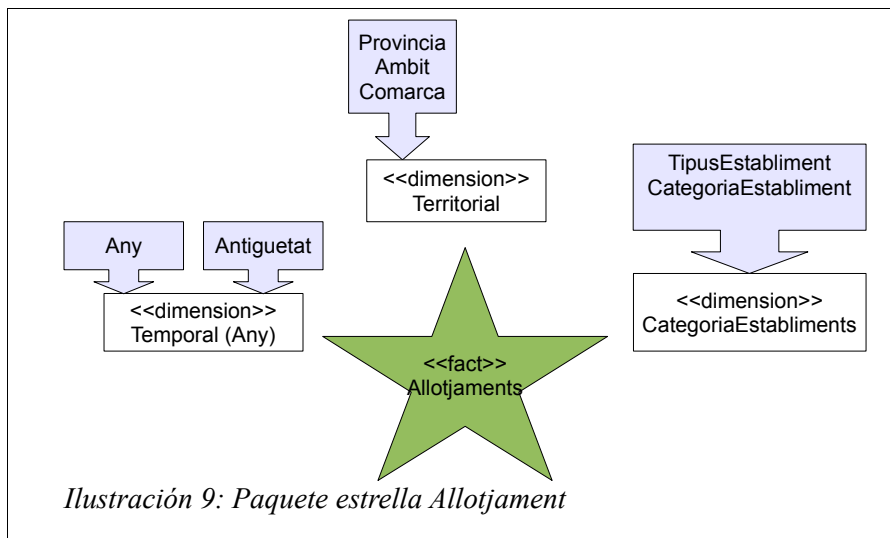
Comarcas y Municipios: La razón de que existan dos puntos desde donde poderse “enganchar” al árbol de dimensiones territorial, es debido a que existen, *hechos* que tienen su nivel más bajo de

representación en el *Municipi* y otros en *Comarca*.

En el anexo (Anexo Modelo Conceptual basado en UML - Nivel 2: Paquete dimensión Comuns) del detalle en UML existe un modelado detallado del paquete

### 2.4.5. NIVEL 2: PAQUETE ESTRELLA ALLOTJAMENTS

La estrella *Allotjaments* tiene como el *hecho* de que exista un establecimiento y su número de plazas. No se puede acumular dichos valores por varios años ya que no tiene sentido, en todo caso establecer valores medios.



Las dimensiones tienen como jerarquía de base, el año (*Any*), comarca (*Comarca, Ambit, Provincia*) y categoría de establecimiento (*CategoriaEstabliment: TipusEstabliment, CategoriaEstabliment*)

La dimensión *CategoriaEstabliments* tiene la siguiente estructura.

<i>TipusEstabliment</i>	<i>CategoriaEstabliment</i>
<i>Hotels</i>	<i>Hotels (estrellas or)</i>
	<i>Hostals (estrellas argent)</i>
<i>Campings</i>	<i>Luxe</i>
	<i>1</i>
	<i>2</i>
	<i>3</i>
<i>Turisme Rural</i>	<i>Casa de poble compartida.</i>
	<i>Casa de poble independent</i>
	<i>Masia</i>
	<i>Masoveria</i>

*Tabla 2: Tabla de tipos y categorías de Establecimiento*

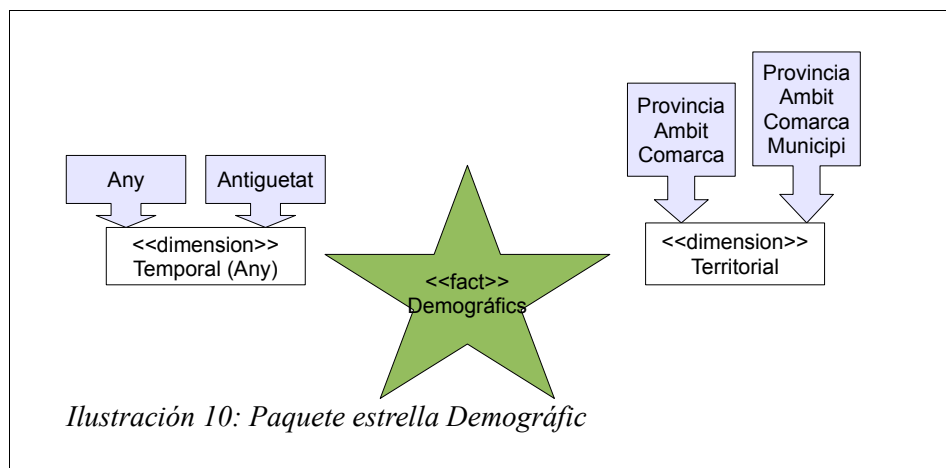
En el anexo (Anexo Modelo Conceptual basado en UML - Nivel 2: Paquete estrella



Allotjaments) del detalle en UML existe el correspondiente modelado detallado del paquete

#### 2.4.6. NIVEL 2: PAQUETE ESTRELLA DEMOGRÁFICS

La estrella *Demogràfics* tiene como el hecho a los datos demográficos de una zona geográfica (*habitants, homes, dones*) y, por otra parte, el dato geográfico de los Km2 cuadrados (*extensio.*) En el caso de *homes* o *dones* las fuentes de datos solo suministran información del año 2012, pero el sistema está preparado para insertarla en otros años.



Tenemos a las dimensiones temporal (*any*) y territorial sobre las que se puede analizar los hechos. En el caso del nivel territorial, puede ser tanto a nivel municipio (*municipi, comarca, ambit, província*) como a nivel comarca (*comarca, ambit, província*). En los informes las correlaciones *demogràfics* con *establiments* se realizará a nivel *comarca*

No se puede acumular los valores demográficos y geográficos por varios años ya que no tiene sentido, si que se puede calcular las medias, aunque en el caso de los Km<sup>2</sup> de un municipio no sea un dato que experimente cambio.

En el anexo ( Anexo Modelo Conceptual basado en UML - Nivel 2: Paquete estrella Demogràfics) del detalle en UML existe el correspondiente modelado detallado del paquete

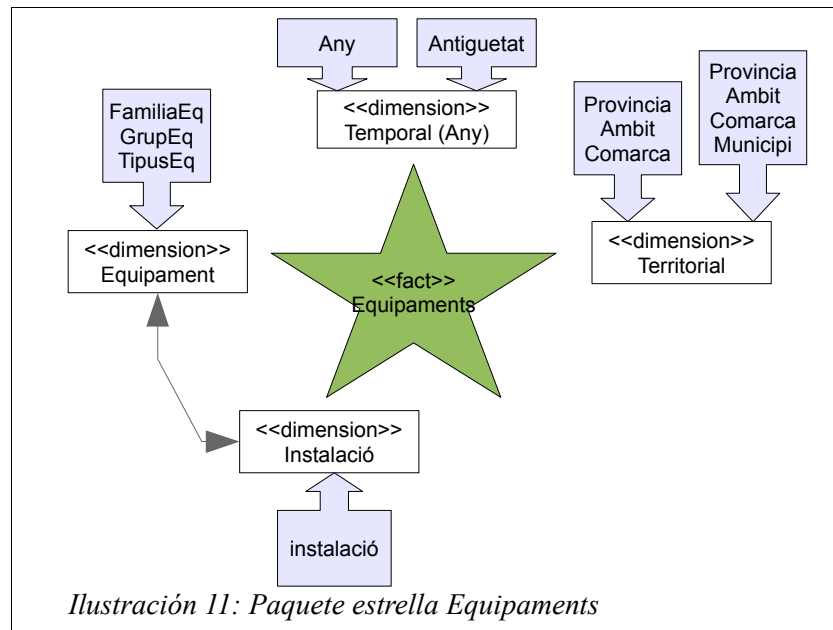
#### 2.4.7. NIVEL 2: PAQUETE ESTRELLA EQUIPAMENTS

La estrella *Equipaments* tiene como hecho la existencia de un determinado tipo de equipamiento y del tipo de instalaciones que posee. No se puede acumular los valores los equipamientos por varios años ya que no tiene sentido. Este es un “hecho” que tiene datos “degenerados de dimensión” (*nom, cp, telefon, longitud, latitud*), pero el valor es la cantidad de equipamientos que existen en una determinado tipo de equipamiento.

En las fuentes de datos la información esta reducida a un solo año, pero el sistema se a diseñado para que permita la entrada de otros años por eso se ha incluido la dimensión temporal (*Any*). De echo se cree, que enriquecería los análisis correlativos entre equipamientos y demográficos/establecimientos, si en un futuro, se inserta la lista de equipamientos de otros años.

Una observación a pesar de que incluya la dimensión temporal, los informes actuales no lo tiene en cuenta, o dicho de otro modo, cualquier correlación con demográficos o establecimientos en

donde intervenga la dimensión temporal, *Equipaments* tendrá los datos de 2012



Un equipamiento está relacionado con la dimensión territorial, tanto a nivel municipio (*municipi, comarca, ambit, província*) como a nivel comarca (*comarca, ambit, província*). En los informes, las correlaciones *equipaments* con *establiments* se realizará a nivel *comarca*. También es a nivel *comarca* en el caso de las correlaciones *equipaments* con *demogràfics*, aunque ambos posean el nivel municipio, ya que *equipaments* tiene muchos más municipios que *demogràfics*.

El análisis de las fuentes de datos nos da que la estructura de clasificación de los equipamientos (*equipament <dimension>*) está formada por 3 niveles jerárquicos que los hemos identificado como *FamiliaEq, GrupEq* y *TipusEq*.

A mayores tras estos 3 niveles y en función de *TipusEq* puede el equipamiento contener unas instalaciones (*instalacions*) para realizar una serie de actividades. El concepto de instalaciones puede hacer referencia tanto a los espacios como a los servicios. Un ejemplo del segundo caso son los centros educativos que pueden dar unos o otros niveles de enseñanza.

La relación entre la dimensión de equipamientos (*equipament <dimension>*: *FamiliaEq, GrupEq* y *TipusEq*) e instalaciones es una relación N a M, es decir, una instalación puede estar en múltiples *equipament <dimension>* y a su vez *equipament <dimension>* puede tener múltiples instalaciones.

En el anexo (Anexo Modelo Conceptual basado en UML - Nivel 2: Paquete estrella Equipaments) del detalle en UML existe el correspondiente modelado detallado del paquete.

### 3. Capítulo. Diseño

#### 3.1. MODELO LÓGICO MULTIDIMENSIONAL.

El modelo lógico multidimensional se diseña siguiendo lo que sería el modelo lógico en el diseño de una base de datos relacional.

### 3.1.1. CRITERIOS GENERALES PARA EL MODELO LÓGICO MULTIDIMENSIONAL.

Como se ha comentado en el Modelo Conceptual el proceso de definir el modelo lógico multidimensional se ha realizado con la misma herramienta para el modelado conceptual **ArgoUML**, ya que como se explica previamente, uno de los objetivos del proyecto era tener un método y un modelado unificado e integrado.

En sección anexa (Anexo Modelo Lógico basado en UML) hay una explicación detallada del como se hizo el modelado, de como se explicita la relación entre los elementos del conceptual con los elementos del lógico, lo que posibilita la integridad entre ambos modelos. Además se creó un perfil de estereotipos específicos para poder representar tanto el modelo como su relación con el modelo conceptual.

### 3.1.2. NOTAS SOBRE LA REPRESENTACIÓN DEL MODELO LÓGICO

La representación del modelo lógico se ha extraído del repositorio UML,(ver anexo Anexo Modelo Lógico basado en UML) pero es una versión preparada para la exposición de la memoria. En ella utiliza muchos de los artefactos del perfil diseñado para poder representar el modelo de base de datos. Así tenemos:

- `<<tabla>>` a nivel de clase para indicar que es un tabla.
- `<<fk>>` a nivel atributo indica que es un clave foranea.
- `<<pk>>` a nivel atributo indica que es un clave primaria
- `<<unsigned>>` hace referencia a que es tipo numérico sin signo.

### 3.1.3. NOMENCLATURA DE LOS ELEMENTOS

Definimos en el modelo lógico los nombres finales que poseerán los diversas tablas y campos, por lo que se han establecido una serie de convenciones, siguiendo un versión libre del modelo propuesto en [ref-Bouman] .

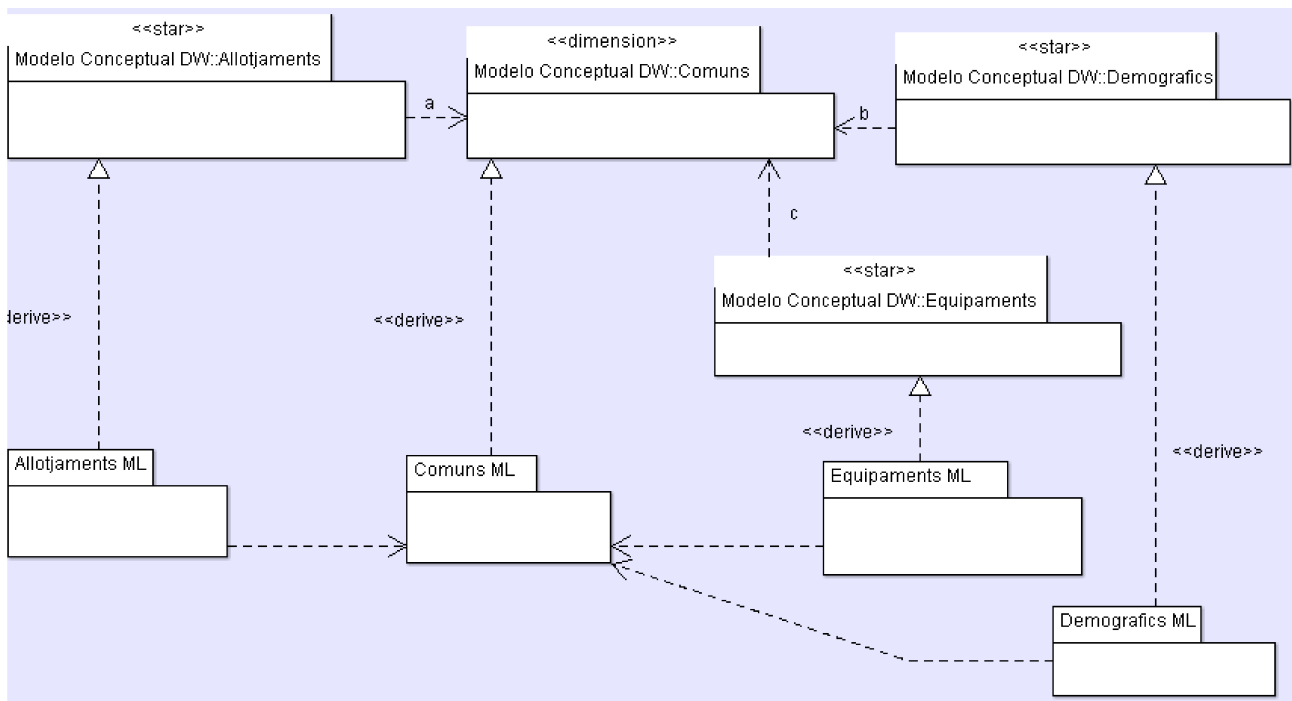
- prefijo **dim** indica que la tabla es una dimensión.
- prefijo **fact** indica que la tabla es un hecho o valor
- prefijo **agg** indica que es tabla con valores agregados de otras tablas.
- prefijo **deg** indica que es un tabla resultado de una “*dimensión degenerada*”
- prefijo **clau** en un atributo indica que es una columna que forma parte de la clave.

### 3.1.4. MODELO LÓGICO NIVEL 1: VISIÓN GENERAL.

En el diagrama de la visión general del modelo lógico del almacén de datos, nos permite ver la derivación (relación) que existe con el modelo conceptual. Para buscar ser homogéneos se han aplicado los mismos criterios en la división en paquetes. Con ello, además, facilitamos el mostrar la relación entre los modelos. Así tenemos:

- Del modelo conceptual *Comuns* deriva el modelo lógico *Comuns*.
- Del modelo conceptual *Allotjaments* deriva el modelo lógico *Allotjaments*

- Del modelo conceptual *Demogràfic* deriva el modelo lógico *Demogràfic*
- Del modelo conceptual *Equipaments* deriva el modelo lógico *Equipaments*



Il·lustració 12: Nivel 1 Modelo lógico

### 3.1.5. NIVEL 2: PAQUETE LÓGICO COMÚNS.

En la diseño de la tabla del año (*dim\_any*) se ha incluido el nuevo atributo *actual* que sirve para indicar cual es el año actual y cuales son los años anteriores, esto, por ejemplo, aporta rapidez en las selecciones del tipo comparativa entre este año y los 3 anteriores.

En el diseño de la tabla de comarca (*dim\_comarca*) se ha decido desnormalizar y agrupar todos los elementos en un sola tabla y que en ella estén todas las combinaciones posibles. En esté caso al ser jerárquico solo hay una combinación por *comarca*, *ambit* y *provincia*.

Lo mismo pasa con la tabla *dim\_muicipi* que incorpora *comarca*, *ambit* y *provincia*. También se ha decido que una fila de *dim\_municipi* solo puede existir si existe el correspondiente fila en *dim\_comarca*.

Las claves de *dim\_comaca* y *dim\_muicipi* serán autogeneradas por el sistema gestor de bases de datos.

La representación de las tablas aparecen en las ilustraciones que acompañan a las descripciones de los demás modelos.

En el anexo (Anexo Modelo Lógico basado en UML - Nivel 2: Paquete Lógico *comúns*) del detalle en UML existe el correspondiente modelado detallado del paquete

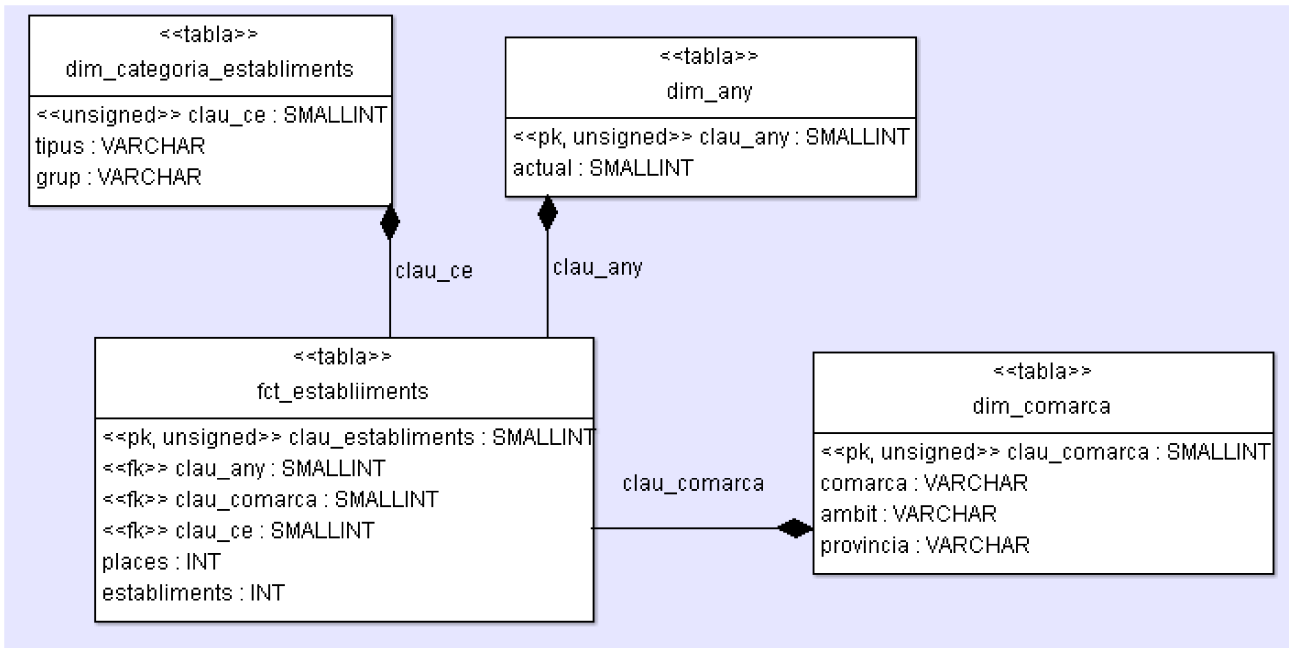
### 3.1.6. NIVEL 2: PAQUETE LÓGICO ALLOTJAMENTS.

En el diseño de la tabla de categorías (*dim\_categoria\_establiments*) se ha decido desnormalizar y agrupar todos los elementos en un sola tabla y que en ella estén todas las combinaciones posibles. En esté caso al ser jerárquico solo hay una combinación por *tipus* y *grup*.

En *fct\_establiments* están los valores del número de plazas (*places*) y el número de establecimientos (*establiments*)

Todas las claves primarias de las tablas sobre alojamientos serán generadas por el sistema gestor de base de datos.

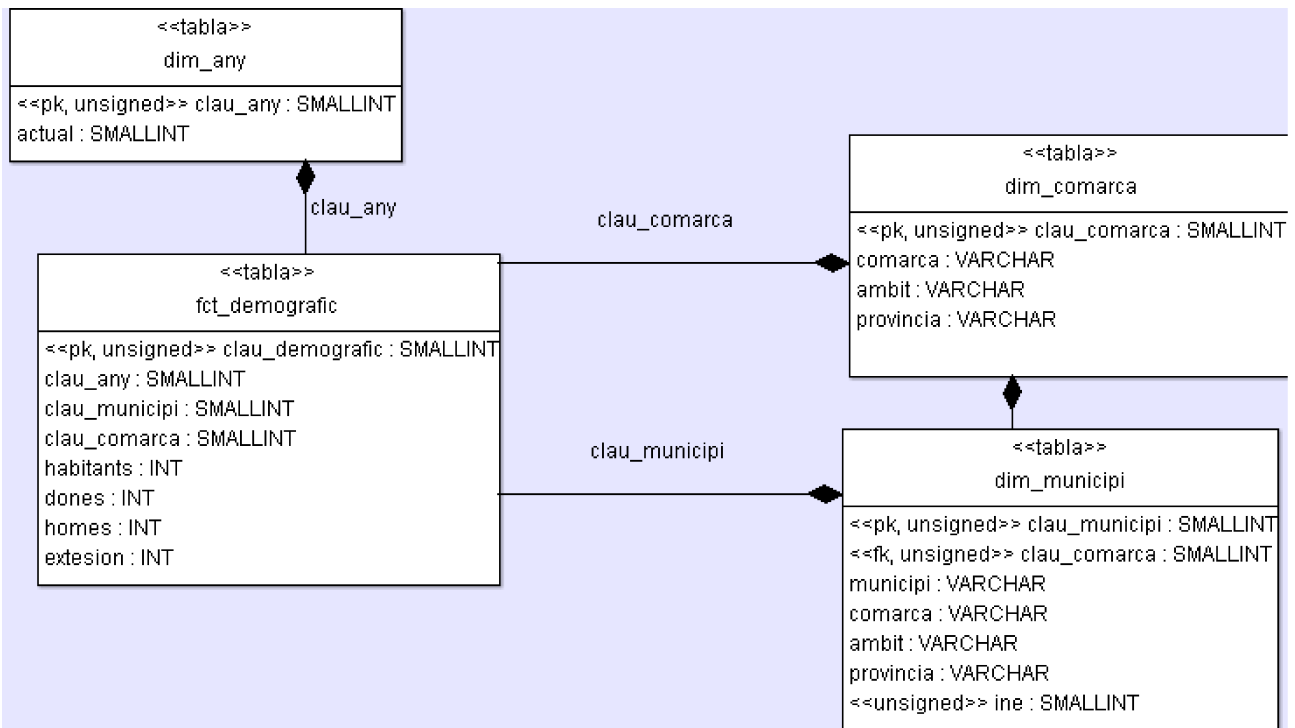
La existencia de un valor en *fct\_establiments* viene dada por la existencia en *dim\_categoria\_establiments*, *dim-any* y *dim\_comarca* de referencia



Il·lustració 13: Paquete lògic Allotjaments

En el anexo (Anexo Modelo Lógico basado en UML - Nivel 2: Paquete Lógico Allotjaments) del detalle en UML existe el correspondiente modelado detallado del paquete

### 3.1.7. NIVEL 2: PAQUETE LÓGICO *DEMOGRAFIC*.



Il·lustración 14: Paquete lógico Demográfico

La clave primarias *fct\_demografic* será generada por el sistema gestor de base de datos.

En *fct\_demografic* están los valores sobre el número de *habitants*, *dones* y *homes* y la *extensió*. Los atributos *Dones* y *homes* pueden tener valores nulos si no se han importado o calculado dichos datos. Ver “[Suposiciones aplicadas](#)” en el “Procesos de extracción, transformación y carga (ETL)”

La existencia de un valor en *fct\_demografic* viene dada por la existencia en *din\_any*, *dim\_municipi*, *dim\_comarca* de referencia

En el anexo (Anexo Modelo Lógico basado en UML - Nivel 2: Paquete Lógico Demografic) del detalle en UML existe el correspondiente modelado detallado del paquete

### 3.1.8. NIVEL 2: PAQUETE LÓGICO *EQUIPAMENTS*.

Como lo que es de información relevante es el número de establecimientos de un tipo y sus instalaciones, se han definido tablas agregadas que contienen dicho valor ya totalizado al nivel más bajo de la jerarquía.



los diferentes elementos que intervienen en dicho proceso. Dicho modelo se ha basado en las propuestas de [ref-LujanS] más las recomendación para un buen proceso de **ETL** de [ref-Bouman]

Se ha decidido crear la definición del proceso **ETL** utilizando el modelado **UML**, con ello se consigue que en un único repositorio de información estén tanto la información sobre el almacén de datos como los procesos de carga y lo que es más importante como se relacionan entre ellos.

Otro aspecto no menos importante es que se definen y se tipifican todos los elementos implicados en un **ETL**, con lo cual, se consigue por un lado normalizar la nomenclatura utilizada y por otra facilitar el análisis y diseño ya que esté se centra en ver como se relacionan y se utilizan dichos componentes para que una fuente de datos concreta se traslade al almacén de datos definido.

El paradigma que persigue el modelo se puede resumir en los siguientes puntos (basados en [ref-LujanS] y [ref-Bouman])

- El proceso será descompuesto en unidades elementales llamados **mecanismos ETL**
- Los datos de cualquier fuente que no estén en tablas será previamente traspasado a una estructura de base de datos relacional.
- Existirá una zona en la que se almacenarán todas las tablas necesarias en el proceso de transformación. *Stagging Area*
- El proceso de transformación, será auditado y en el caso de errores en los datos o en su formato el flujo se detiene en el estado en que se ha producido el error. En nuestro caso cuando se detecte un error se mirará si es factible modificarlo definiendo un paso de conversión en el proceso de transformación.
- Para los fuentes de orígenes no relacionales se realizará un primer paso que será tener dichas fuente en un formato relacional y sin ningún tipo de manipulación. Es importante que los datos y su estructura no sean manipulados, eso se deja a los **mecanismos ETL** de transformación.
- Las tablas intermedias generadas durante el proceso de transformación se inicializarán en cada proceso de importación.

A continuación desglosamos las unidades elementales llamadas **mecanismos ETL** definidos por [ref-LujanS]:

- **Agregación (Aggregation)**, como su nombre indica se refiere cuando se totalizan o calculan datos basados en algún criterio
- **Conversión (Conversion)**: Cambia los tipos de datos y los formatea o genera nuevos datos a partir de datos existentes.
- **Filtro (Filter)**: filtra y verifica los datos.
- **Incorrecto (Incorrect)**: Redirige (almacena) datos incorrectos.
- **Union (Join)** Se unen dos o mas fuentes en una nueva que sea la combinación de algunos de sus atributos.
- **Cargador (loader)** Realiza una carga de datos al almacén de datos.
- **Registrador (log)** Registra la actividad de los demás mecanismos.
- **Fusionar (Merge)** Integra dos o mas fuentes de datos con atributos compatibles



- **Sustituto (*Surrogate*)** Genera una clave primaria sustituta a la que sería la clave natural que define a la entidad. Por ejemplo agrupar los atributos en una clave natural en un única clave en formato numérico.
- **Envoltorio (*wrapper*)** Extrae los datos una fuente de datos externa para insertarlo en los registros de la base de datos del proceso de **ETL**

En el anexo (Anexo Modelo Lógico ETL basado en UML ) existe el correspondiente ampliación de como se utilizó el UML.

### 3.2.1. SUPOSICIONES GENERALES APLICADAS

- Como criterio de diseño se ha establecido que siempre que sea posible utilizar la definición de nombres y las clasificaciones que viene con las fuentes del año 2012. Existe algunos valores que su nombre varía con los años.
- El cálculo de las plazas de camping serán según la regla del 2011, esto implica que en años anteriores durante el proceso de transformación habrá que recalcular las plazas de campings
- El ajustes de los casas rurales en el 2006 ya que no existía la categoría de “*t.rural independent*”
- La inclusión de todos los equipamientos independientemente de que los municipios no tuvieran su correlación con la fuente población.
- El cálculo de *homes y dones* en los años anteriores del 2012 a partir de la proporción del 2012

### 3.2.2. PAQUETE *ETL COMUNS*.

La carga de la tabla *dim\_Any* se produce cuando se decide importar un nuevo fichero csv de un año en concreto. Se indica que es un proceso manual, ya que se entraran los datos del *loader\_any* manualmente para después realizar el proceso de carga, que puede implicar el cambio de los valores del año actual.

La lista de comarcas, ámbitos y provincias se entra manualmente, ya que son pocos datos

Para determinar la lista con los nombres normalizados del municipio se basará en la información que suministran las fuentes de *poblacio.csv* y *equipaments.csv* a las que se *extraerá* en cada uno de ellos el dato de población, para después realizar una *unión* en donde se relaciona el nombre de una fuente (*poblacio.csv*) con el nombre de la otra (*equipaments.csv*) a continuación se realizará un *conversion* para generar el nombre normalizado que será utilizado en el almacén de datos. Con ello tendremos relacionados el nombre normalizado con el nombre de *equipaments.csv* y con el nombre *poblacion.csv*. Durante el proceso de *carga* a *dim\_municipi* se incorporaran los datos de comarca.

En el anexo (Anexo Modelo Lógico ETL basado en UML -Nivel 2: Paquete *ETL Comuns*.) del detalle en UML existe el correspondiente modelado detallado del paquete

### 3.2.3. PAQUETE *ETL ALLOTJAMENTS*.

Es un conjunto de fuentes de datos, uno por cada año, y con estructuras diferentes. Así tenemos un tipo de *envoltorio* para el año 2006, otro tipo para los años 2007,2008,2009,2011 y el año 2012

necesita antes aplicar múltiples transformaciones.

El 2006 necesita una paso adicional de *conversión* al formato de los otros años.

En el caso del 2012 se ha dividido el fichero original en otros 3 ficheros en que se guardan cada uno de los tres grandes bloques que contiene:

- registro de hoteles y campings,
- registro de las casas rurales
- un segundo registro de casas rurales según la asociación de turismo rural

A cada uno ha de tener su propia *conversión* y tras ello se aplica un *fusionar* los datos de las dos fuentes de turismo rural y finalmente un *unión* que lo combine todo y se pueda seguir con el proceso de transformación ya unificado con los otros años.

Una vez unificados los formatos de los distintos años, se aplica un *filtro* para detectar errores, calcular (si es necesario) e insertar listas de ajustes, finalmente se realizará un *transposición* para pasar de columnas por tipo de establecimiento a un fila por cada tipo de establecimiento.

Las datos de sobre las categorías de alojamientos son entrados manualmente y surgen a partir de las cabeceras del año 2012.

En el anexo (Anexo Modelo Lógico ETL basado en UML -Nivel 2: Paquete *ETL Allotjaments*.) del detalle en UML existe el correspondiente modelado detallado del paquete

#### 3.2.4. PAQUETE *ETL DEMOGRAFICS*.

Tras importar *poblacio.csv* se realizará una *transposición* para pasar de los datos de habitantes en columna por años a una fila por año, después se realizará un *filtro* para verificar errores y aplicar ajustes si son necesarios.

Indicar que cuando venga futuros ficheros habrá que revisar si la fuente vuelven a repetir los últimos 6 años o viene solo el nuevo año.

En el anexo (Anexo Modelo Lógico ETL basado en UML -Nivel 2: Paquete *ETL 3.2.4. Demografics*.) del detalle en UML existe el correspondiente modelado detallado del paquete

#### 3.2.5. PAQUETE *ETL EQUIPAMENTS*.

El paquete *ETL Equipaments* se desglosa en dos grandes áreas.. Así tenemos por un lado la carga de los equipamientos y por otro la carga de las categorías de los equipamientos. Primero hay que realizar la importación de las categorías y después de los equipamientos.

Para el proceso de las categorías lo primero es *envolver* la entrada categorías de tal forma que el primero es la *familia* el segundo el *grup* y el tercero es *tipus* y el resto asignarlo a *instalacion\_1... instalacion\_* después hay que seguir dos direcciones:

- *Filtrar* categorías eliminando duplicados y errores. para cargarlo en el *dim\_tipus\_eq*
- *Transponer* instalaciones de columnas a una fila por instalación y *filtrar* eliminando duplicados para *cargarlo* en el *dim\_instalacio*.

Para el proceso de equipamientos, para cada fila hacemos una *unión* con las categorías de equipamientos, aplicamos una *conversión* para colocar el nombre del municipio normalizado, aplicamos *filtros* para detectar duplicados y fallos y finalmente es *cargado* en *fct\_equipment*,

*deg\_equipament* y en *agg\_equipament*

Con los equipamientos correctos e insertados hay que generar la relación entre el equipamiento y las instalaciones, para ello hay que hacer una *unión* incluyendo a la fila todas las claves foráneas y la clave del establecimiento y con la nueva fila obtenida hay que hacer una *transposición* de columnas de instalación a filas para finalmente cargarlo en *fet\_instalacio* y *agg\_instalacio*.

En el anexo (Anexo Modelo Lógico ETL basado en UML -Nivel 2: Paquete *ETL Equipaments*.) del detalle en UML existe el correspondiente modelado detallado del paquete

### 3.3. MODELO FÍSICO.

#### 3.3.1. SISTEMA GESTOR DE BASE DE DATOS.

La dimensión de las tablas y el volumen de los datos hace que no haya problemas de rendimiento con las consultas a las bases de datos.

La descripción de las instrucciones SQL para la generación de tablas del sistema dimensional por falta de tiempo no se han insertado por el momento en este documento. En todo caso seguirán las siguientes consideraciones:

- Como motor de la base de datos se puede utilizar tanto el **MyISAM** como el **InnoDB** este último ofrece almacenamiento transaccional **ACID** y las restricciones de **FOREIGN KEY** pero ambos temas no son fundamentales para un almacén de datos. En cuanto a MyISAM permite la consulta de datos sin procesos de bloqueo y búsquedas a texto completo. En un principio con esta última es suficiente.
- Los índices: Se crearán tantos índices como atributos filtrables existan en las dimensiones y los hechos.
- Como es una base de datos pequeña, no hace falta sistemas de particionado de tablas u otros para optimizar las consultas.
- En los casos de nombres y de descripciones estará establecido tamaño del varchar al máximo. 255, para evitar problemas futuros de truncamiento.

A continuación se detalla la definición de las tablas del almacén de datos, están por orden alfabético:

```
CREATE TABLE `agg_equipament` (  
  `clau_equipament_calcul` smallint(5) unsigned NOT NULL  
  AUTO_INCREMENT,  
  `clau_any` smallint(4) unsigned NOT NULL,  
  `clau_municipi` smallint(5) unsigned NOT NULL,  
  `clau_comarca` smallint(5) unsigned NOT NULL,  
  `clau_tipus_eq` smallint(5) unsigned NOT NULL,  
  `quantitat` int(5) unsigned DEFAULT '0',  
  PRIMARY KEY (`clau_equipament_calcul`),  
  KEY `ibuscac` (`clau_any`, `clau_municipi`, `clau_comarca`,  
  `clau_tipus_eq`)  
) ENGINE=InnoDB DEFAULT CHARSET=utf8;
```

```
CREATE TABLE `dim_any` (  
  `clau_any` smallint(4) unsigned NOT NULL,  
  `actual` smallint(2) NOT NULL DEFAULT '0',  
  PRIMARY KEY (`clau_any`)  
) ENGINE=InnoDB DEFAULT CHARSET=utf8;
```

```
CREATE TABLE `dim_comarca` (  

```

```
CREATE TABLE `agg_instalacio` (  
  `clau_equipament_calcul` smallint(5) unsigned NOT NULL,  
  `clau_instalacio` smallint(5) unsigned NOT NULL,  
  `clau_any` smallint(4) unsigned NOT NULL,  
  `clau_municipi` smallint(5) unsigned NOT NULL,  
  `clau_comarca` smallint(5) unsigned NOT NULL,  
  `clau_tipus_eq` smallint(5) unsigned NOT NULL,  
  `quantitat` int(5) unsigned DEFAULT '0',  
  PRIMARY KEY (`clau_equipament_calcul`, `clau_instalacio`)  
) ENGINE=InnoDB DEFAULT CHARSET=utf8;
```

```
CREATE TABLE `dim_categoria_establiments` (  
  `clau_ce` smallint(5) unsigned NOT NULL AUTO_INCREMENT,  
  `tipus` varchar(255) NOT NULL,  
  `grup` varchar(255) NOT NULL,  
  PRIMARY KEY (`clau_ce`)  
) ENGINE=InnoDB DEFAULT CHARSET=utf8;
```

```
CREATE TABLE `dim_instalacio` (  

```

```
`clau_comarca` smallint(5) unsigned NOT NULL
AUTO_INCREMENT,
`comarca` varchar(255) NOT NULL,
`ambit` varchar(255) NOT NULL,
`provincia` varchar(255) NOT NULL,
PRIMARY KEY (`clau_comarca`),
KEY `icomarca` (`comarca`),
KEY `iambit` (`ambit`),
KEY `iprovincia` (`provincia`),
KEY `iterritori` (`provincia`,`ambit`,`comarca`)
) ENGINE=InnoDB DEFAULT CHARSET=utf8;
```

```
CREATE TABLE `dim_municipi` (
`clau_municipi` smallint(5) unsigned NOT NULL
AUTO_INCREMENT,
`municipi` varchar(255) NOT NULL,
`ine` smallint(5) unsigned NOT NULL,
`version` datetime DEFAULT NULL,
`comarca` varchar(255) NOT NULL DEFAULT 'no definido',
`ambit` varchar(255) NOT NULL DEFAULT 'no definido',
`provincia` varchar(255) NOT NULL DEFAULT 'no definido',
PRIMARY KEY (`clau_municipi`),
UNIQUE KEY `municipi` (`municipi`),
KEY `iterritori` (`provincia`,`ambit`,`comarca`)
) ENGINE=InnoDB DEFAULT CHARSET=utf8;
```

```
CREATE TABLE `fct_demografic` (
`clau_demografic` smallint(5) unsigned NOT NULL
AUTO_INCREMENT,
`clau_any` smallint(4) unsigned NOT NULL,
`clau_municipi` smallint(5) unsigned NOT NULL,
`clau_comarca` smallint(5) unsigned NOT NULL,
`habitants` int(5) NOT NULL DEFAULT '0',
`dones` int(5) DEFAULT '0',
`homes` int(5) DEFAULT '0',
`extension` int(5) NOT NULL DEFAULT '0',
PRIMARY KEY (`clau_demografic`)
) ENGINE=InnoDB DEFAULT CHARSET=utf8;
```

```
CREATE TABLE `fct_establiments` (
`clau_establiments` smallint(5) unsigned NOT NULL
AUTO_INCREMENT,
`clau_any` smallint(4) unsigned NOT NULL,
`clau_comarca` smallint(5) unsigned NOT NULL,
`clau_ce` smallint(5) unsigned NOT NULL,
`places` int(5) NOT NULL DEFAULT '0',
`establiments` int(5) NOT NULL DEFAULT '0',
PRIMARY KEY (`clau_establiments`)
) ENGINE=InnoDB DEFAULT CHARSET=utf8;
```

```
`clau_instalacio` smallint(5) unsigned NOT NULL
AUTO_INCREMENT,
`instalacio` varchar(255) NOT NULL,
`version` datetime DEFAULT NULL,
PRIMARY KEY (`clau_instalacio`),
KEY `idx_dim_instalacio_lookup` (`instalacio`)
) ENGINE=InnoDB DEFAULT CHARSET=utf8;
```

```
CREATE TABLE `dim_tipus_eq` (
`clau_tipus_eq` smallint(5) unsigned NOT NULL
AUTO_INCREMENT,
`familia` varchar(255) NOT NULL,
`grup` varchar(255) NOT NULL,
`tipus` varchar(255) NOT NULL,
`version` datetime DEFAULT NULL,
PRIMARY KEY (`clau_tipus_eq`),
KEY `idx_dim_tipus_eq_lookup` (`familia`,`grup`,`tipus`)
) ENGINE=InnoDB AUTO_INCREMENT=112 DEFAULT
CHARSET=utf8;
```

```
CREATE TABLE `fct_equipment` (
`clau_equipment` smallint(5) unsigned NOT NULL
AUTO_INCREMENT,
`clau_any` smallint(4) unsigned NOT NULL,
`clau_municipi` smallint(5) unsigned NOT NULL,
`clau_comarca` smallint(5) unsigned NOT NULL,
`clau_tipus_eq` smallint(5) unsigned NOT NULL,
`version` datetime DEFAULT NULL,
`clau_deg_equipment` smallint(5) DEFAULT NULL,
PRIMARY KEY (`clau_equipment`),
UNIQUE KEY `ibuscar`
(`clau_any`,`clau_tipus_eq`,`clau_deg_equipment`),
KEY `idx_fct_equipment_lookup`
(`clau_equipment`,`clau_any`),
KEY `iagregar` (`clau_any`,`clau_municipi`,`clau_comarca`,
`clau_tipus_eq`)
) ENGINE=InnoDB DEFAULT CHARSET=utf8;
```

```
CREATE TABLE `fct_instalacio` (
`clau_equipment` smallint(5) unsigned NOT NULL,
`clau_instalacio` smallint(5) unsigned NOT NULL,
PRIMARY KEY (`clau_equipment`,`clau_instalacio`)
) ENGINE=InnoDB DEFAULT CHARSET=utf8;
```

La descripción de las tablas del proceso ETL, se aplicará lo expuesto en su modelo lógico creándose en el momento del proceso, pero sus tablas se mantendrán para futuras incorporaciones. .

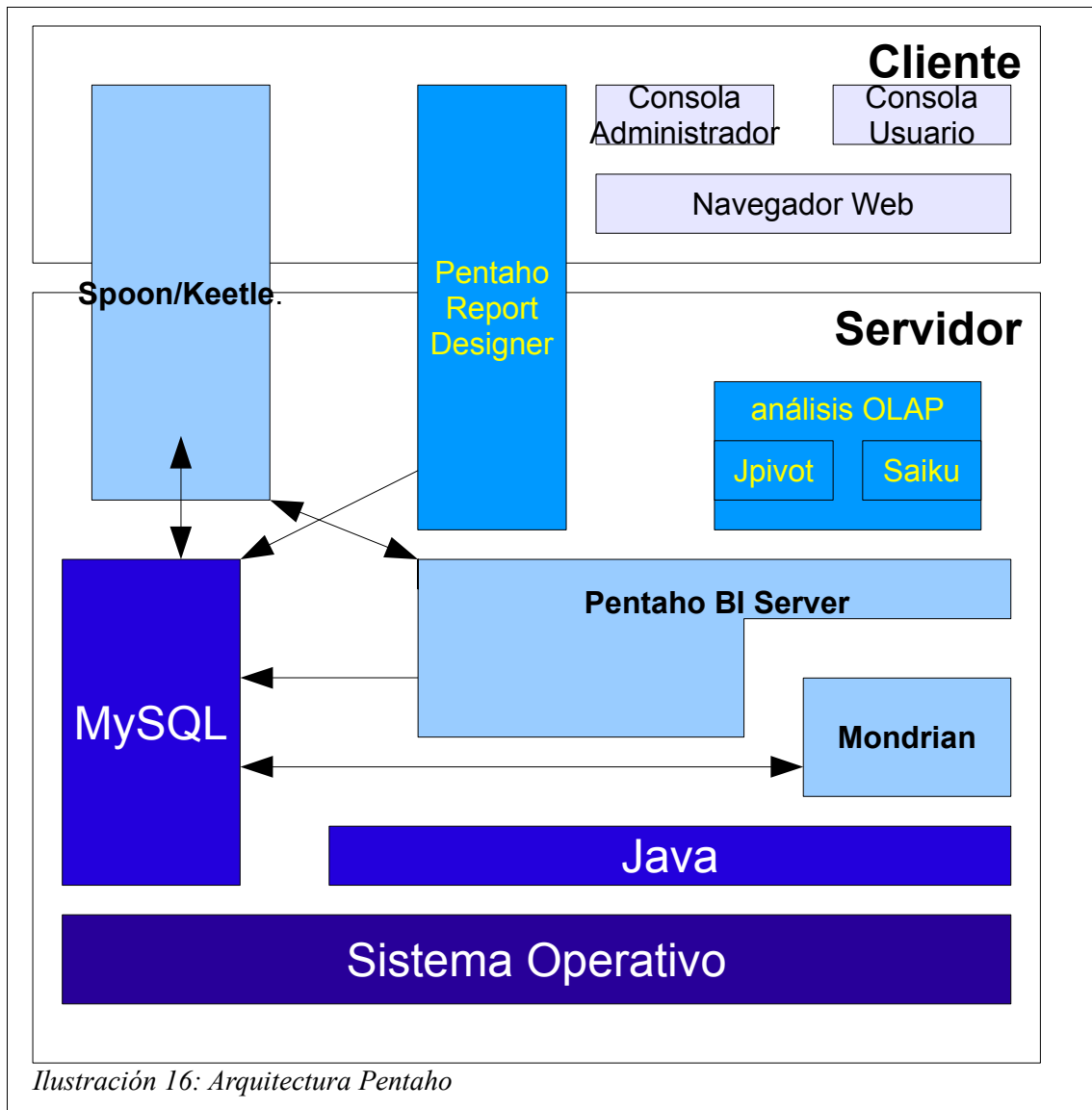
Para claridad y diferenciar existirán dos bases de datos. una para el almacén de datos y otra para el proceso ETL, ambas estarán gestionadas por el mismo servidor de bases de datos, lo que facilitará el proceso de carga y consulta al almacén de datos.

### 3.3.2. EQUIPAMIENTO

Todo el software utilizable está concebido en una arquitectura cliente/servidor. Así que, para este proyecto el equipo servidor puede compartir recursos con otros servicios, ya que, el volumen no es grande y tampoco se espera una carga de trabajo intensa.

Solo en el caso de que se aportaciones de datos con volúmenes mayores e información mas compleja se podría pensar en utilizar un servidor dedicado (virtual o físico)

En el caso de que el producto creado sea un producto de un solo uso se puede pensar en utilizar los sistemas en la nube para alojar los datos y pagar por solo el tiempo de uso.



En cualquier caso la estructura ha de ser de cliente/servidor estando en la parte cliente las herramientas consulta y administración del servidor, ambas se ejecutan sobre cualquier navegador que permita AJAX, se podrán utilizar las herramienta de ETL Spoon / Keettle tanto en un equipo cliente, como en el mismo equipo servidor, lo mismo se puede decidir para Pentaho ReportDesigner y tanto el servidor de Pentaho BI como el de las base de datos ha de estar en un equipo con funciones de servidor (virtual o físico)

Al ser la pila de productos de Pentaho desarrollado en java se puede poner en cualquier S.O. (Windows, Unix, Linux, Mac) que soporte la maquina virtual java.

El MySQL se puede instalar entre otros en equipos Windows y Linux.

## 4. Capítulo Implementación

---

### 4.1. EXPOSICIÓN DEL TRABAJO DE INTEGRACIÓN DE LOS DATOS.

Todo el proceso de transformación se ha realizado en **Spoon** el *Data Integracion de Penthao* en una *Staging Area* en el servidor **MySQL** llamada *tfctl* con los mismos permisos de usuario que las demás bases de datos: `root > tfcmmd`

El primer paso es insertar los ficheros planos a la base de datos es realizado de forma manual, también es manual la separación de los datos de *establiments* del 2012 en 3 ficheros, el volcado a la base de datos se realiza mediante la orden de importación del **MySQL**, teniendo en cuenta los distintos tipos de caracteres que tienen cada uno de ellos.

Ejemplos de los comandos utilizados:

```
load data infile '/tmp/TFC/Equipaments.csv' into table equipaments fields terminated by ',' OPTIONALLY ENCLOSED BY '"' IGNORE 1 LINES;
```

```
load data infile '/tmp/TFC/est2006.csv' into table allotjaments_2006 fields terminated by '\t' TERMINATED BY '\r\n';
```

```
load data infile '/tmp/TFC/poblacioUTF8.csv' into table poblacio fields terminated by ',' IGNORE 1 LINES;
```

En este caso se hizo una transformación primero del fichero UTF8

A partir de ese momento en que están los datos cargados en la base de datos *tfctl* se inicia todo el proceso de transformación con la herramienta **Spoon**.

En **Spoon** se ha creado un trabajo que sirve de guía y orientación de los pasos que hay que hacer, no es recomendable lanzarlo desde ahí (de hecho el siguiente paso a inicio es una cancelación de la tarea, impidiendo realizar los siguientes trabajos) Es mejor lanzar uno a uno, y más recomendable es realizar un seguimiento mediante la ventana de logs "*job metrics*", para ello es mejor lanzar cada transformación desde la misma pagina de transformación. Así se obtienen información detallada de las "métricas del proceso" que muy son útiles ya que informan de los registros leídos, escritos, actualizados o borrados y con ello se tiene cumplida cuenta que se están procesando todos los registros entrados.

Para el proceso se ha seguido lo propuesto el documento lógico ETL del diseño con los ajustes para adaptarlo a la herramienta **Spoon**.

Si entramos en el **Spoon** y revisamos el trabajo etiquetado como "*Panel de trabajos migración TFCDW*" vemos que se ha dividido en diversas tareas y estas a su vez en múltiples transformaciones.

Tememos como grandes bloques: las transformación "*dim\_any*" y "*insertar Categories i Comarques*" y las transformaciones "*treball dim\_poblacio*", "*treball fct\_demografic*", "*treball allotjaments*", "*treball equipaments*".

A continuación se van a detallar dichas transformaciones agrupándolas por las misma división que se hizo en paquetes durante el análisis y el diseño: *comunes*, *allotjaments*, *demografics* y *equipaments*

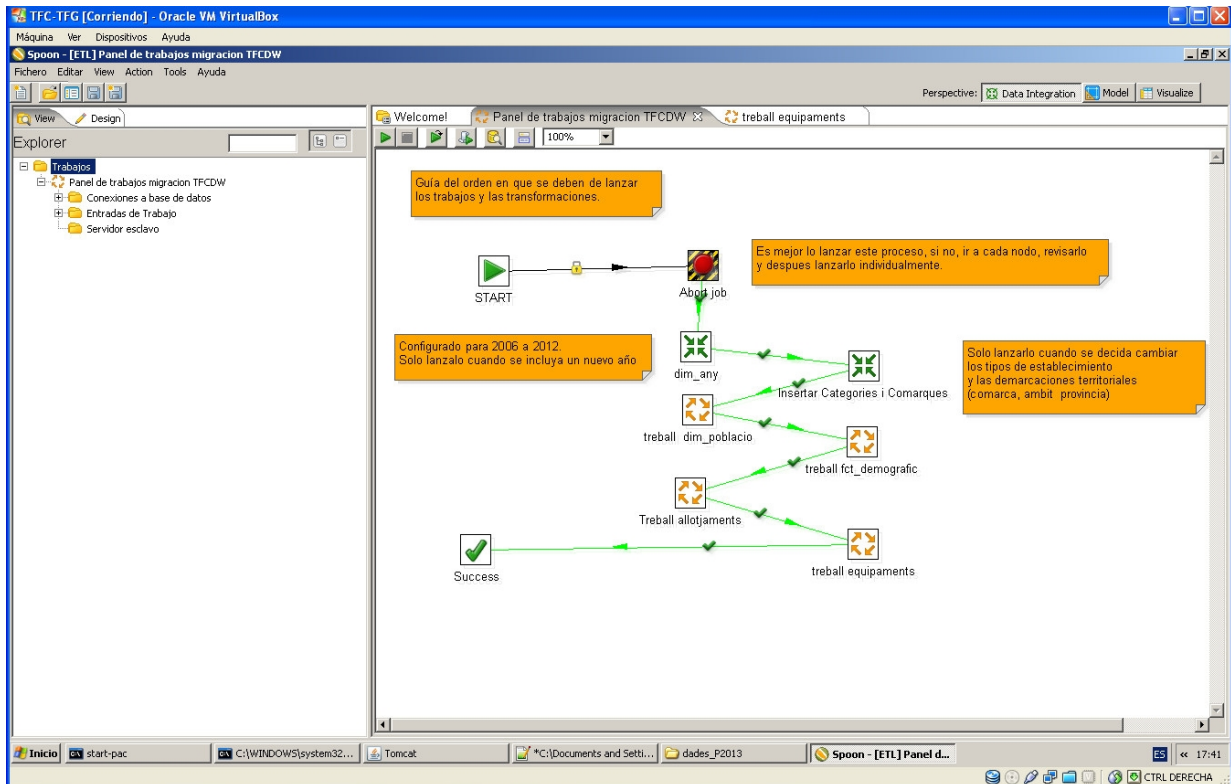


Ilustración 17: Imagen del trabajo que funciona como guía de los pasos a realizar

#### 4.1.1. DATOS COMUNES.

### Transformación para la dimensión temporal (DIM\_any)

En este paso se aprovecha las funcionalidades de generar filas y secuenciar de **Spoon** para crear una secuencia desde el 2006 a 2012 y con actual de 0 a -6.

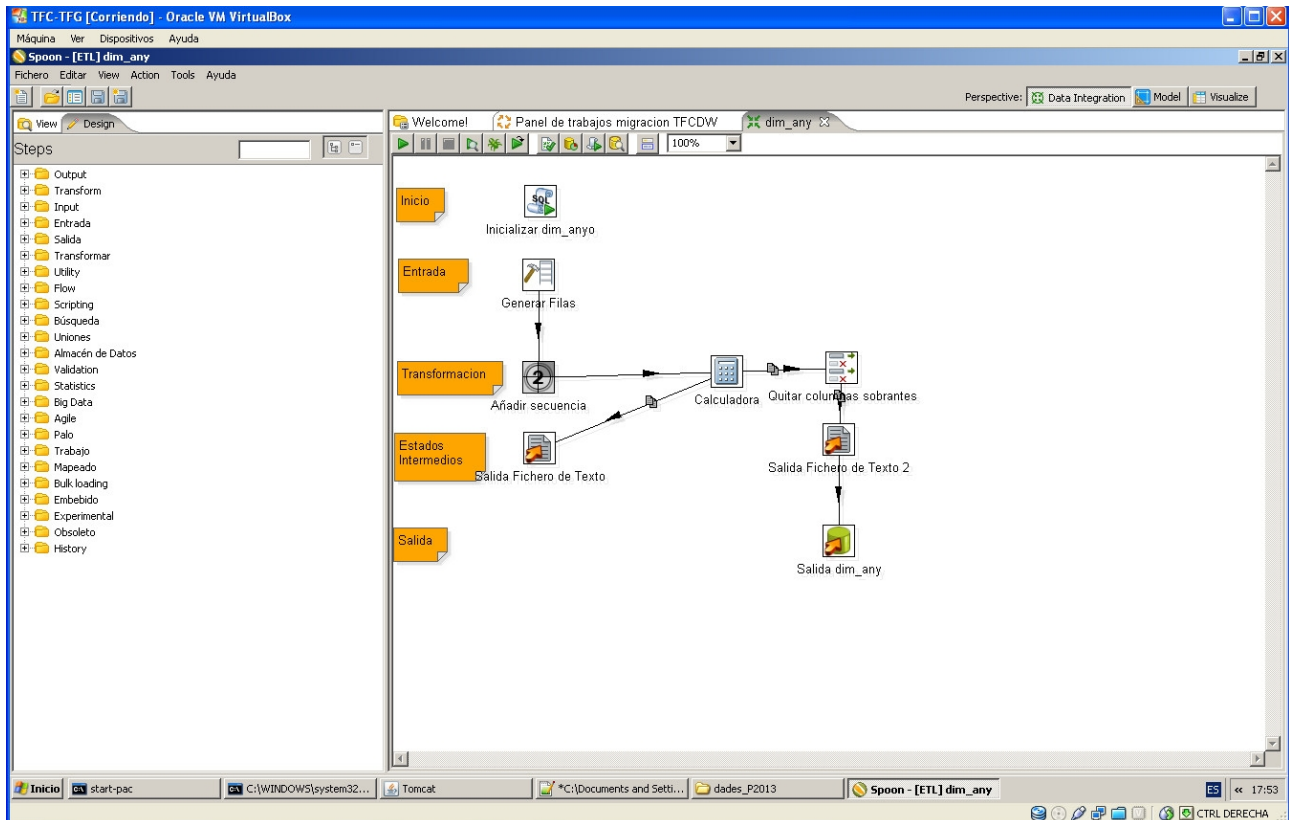


Ilustración 18: Generación de la dimensión temporal

### Transformación para *Categorías de establiments i Comarques*:

Contiene dos scripts SQL que mediante “**inserts**” SQL se entran las tablas de comarca y categorías de establecimientos.

Los nombres de las comarcas se han sacado de la documentación oficial, y su asignación al “*ambit*” se realiza según la convención del 2012.

### Transformación *Municipi dim\_poblacio*:

Responsable de identificar los municipios, tanto si provienen de equipamientos como de población, además tiene la responsabilidad de crear los ficheros intermedios que serán utilizados para recuperar correctamente el municipio durante las transformaciones para los datos demográficos y de equipamientos.



Tal como se indico en el modelo Lógico existe un fichero intermedio que contiene el nombre del municipio y comarca que se "viene" de cada una de las fuentes y unas columnas a mayores con el nombre y comarca "definitivos" que se utilizarán en el sistema.

Por lo tanto, hay que siempre asegurarse que se ha lanzado este paso y de que los ficheros intermedios contienen la información correcta.

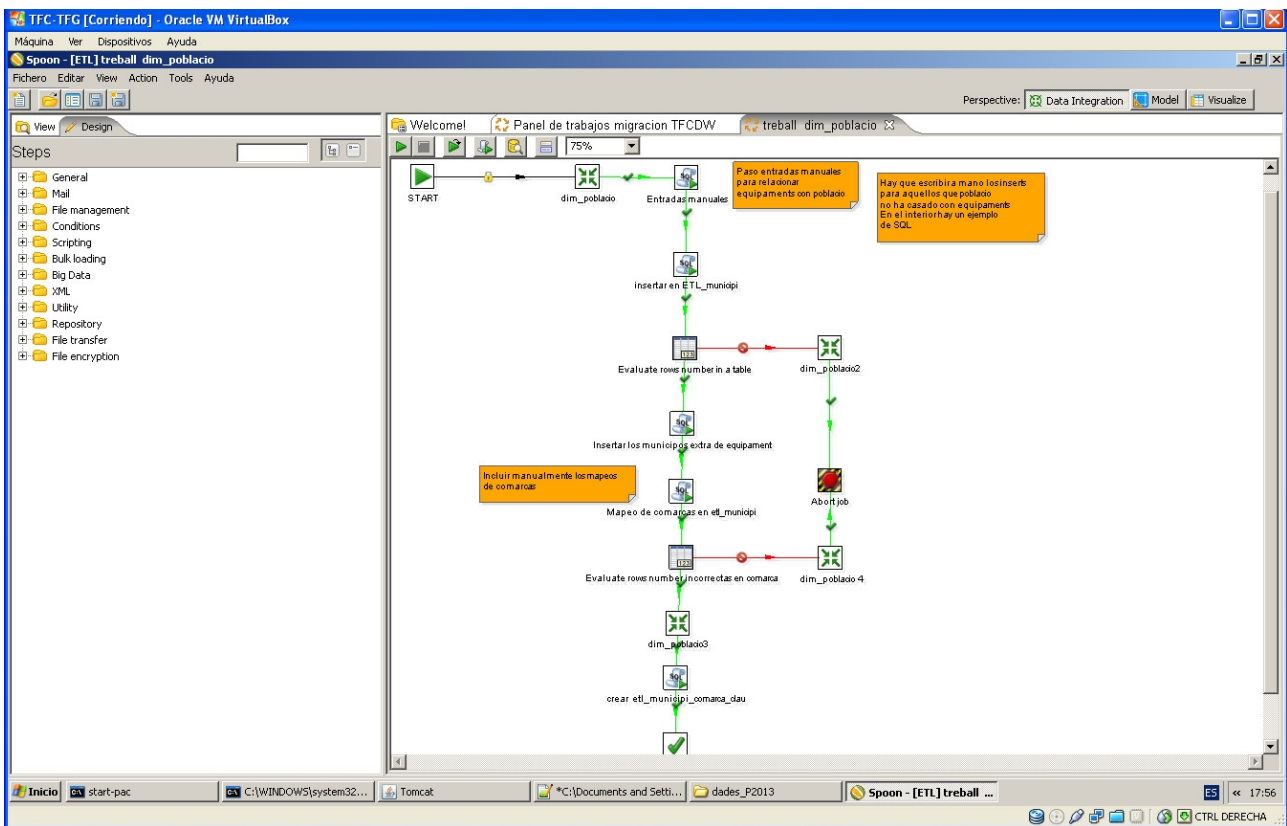


Ilustración 19: Generación de nombres de municipios

Como observación indicar que hay un número de municipios que están en *equipaments* y no en *població*, aún así se han incorporado en el sistema ya que se va a trabajar principalmente con información a nivel de comarca. Pero es importante saber que en habrán más equipamientos que los que hubiese si solo se daban por correctos los municipios existentes en demografía.

#### 4.1.2. DEMOGRÁFICOS.

Se realiza primero una transposición, para después limpiar de puntos y valores incorrectos en los campos de números, por ejemplo: los "nd" se han pasado a 0. Se recupera las claves para comarca y municipio y se verifica que lo tengan todos los registros (si el proceso de población ha ido bien ese paso no debería de dar error).

Finalmente se actualiza el *fact\_demografic* y como un último paso debido a que los datos de los *homes* y *dones* solo están en el 2012 se ejecuta un **SQL script** para recalcular los de los años anteriores a partir de la proporción dada en el 2012

Si *d* se corresponde al registro del año que se esta procesando y *d12* corresponde al registro del año 12 tenemos.

$$d.homes = TRUNCATE(d.habitans * d12.homes / d12.habitans, 0)$$

$d.dones = d.habitants = TRUNCATE(d.habitans * d12.homes / d12.habitants,0)$

Observación el municipio **Canonja** solo tiene datos a partir del 2012

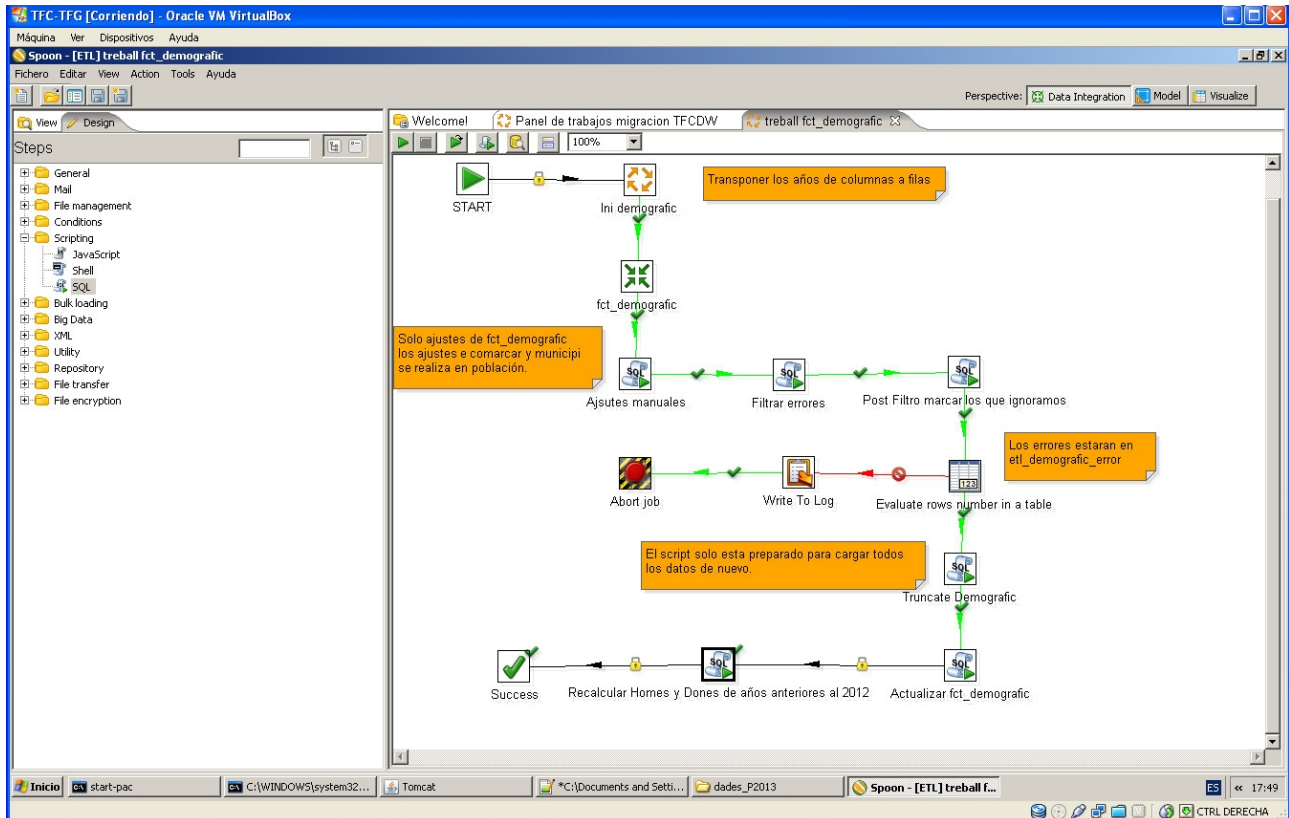
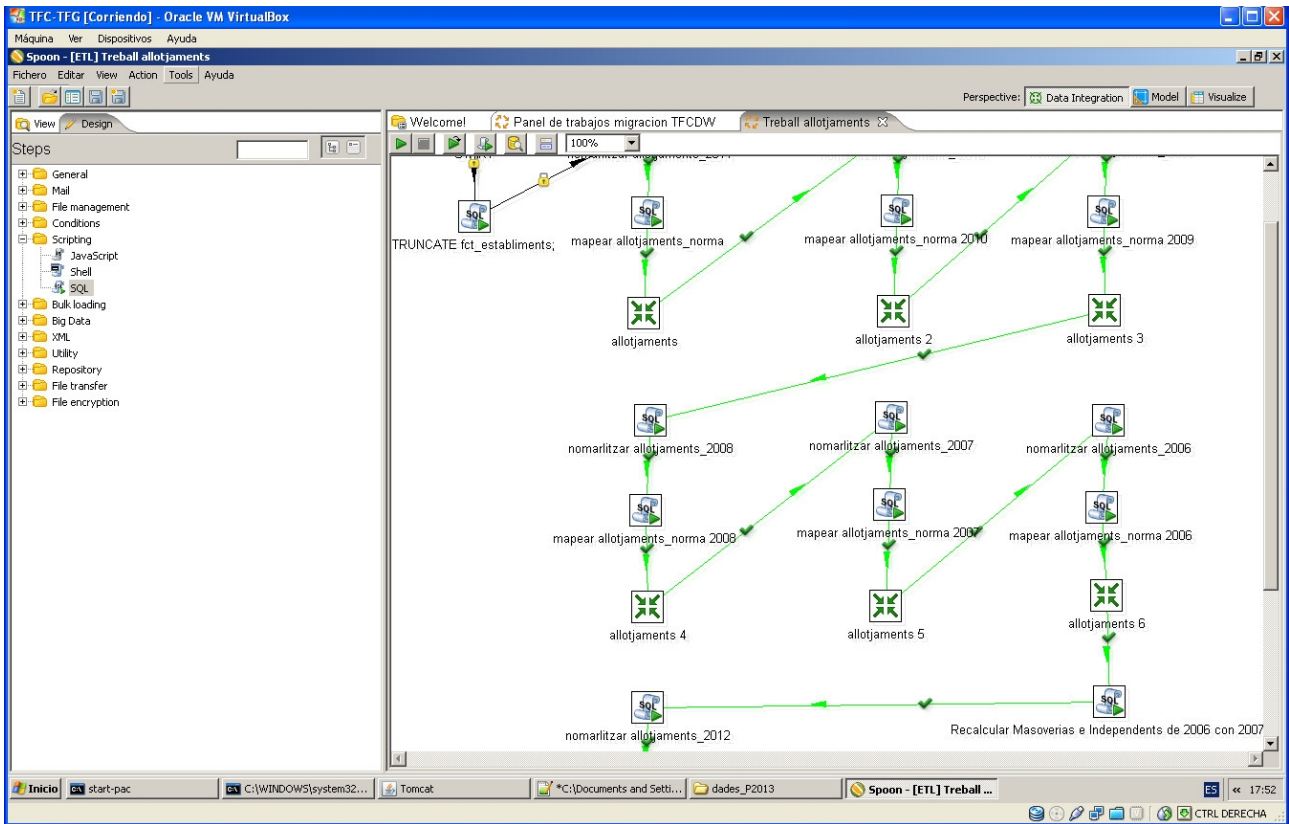


Ilustración 20: Transformación de datos demográficos

### 4.1.3. ALOJAMIENTOS

En el año 2006 los '*campigs privats*' son asignados camping de lujo y se recalcula las *Masoverias e' independents* con los datos del 2007: se pasa el mismo número de *masoveria* del 2007 al 2006 y el resto pasan a ser *independents*.

Para los demás para cada caso se normaliza la entrada de datos, es decir se pasa a una tabla intermedia que contiene las columnas normalizadas y se incluyen los ajustes manuales necesarios mediante scripts SQL. En el caso del 2012 se comparan las dos fuentes de datos que existen para los campings.



Il·lustració 21: Transformació allotjaments

Després cada uno de los años (ver imagen en la siguiente página) se pasa al proceso de filtrado y validación y mediante un *javascript* se validan que los totales dentro del registro cuadren, o al menos que uno de ellos lo haga.

Se valida que los índices a *comarca* sean los correctos y se tipifica el registro de entrada indicando si se corresponde a una comarca, a un ámbito, a una provincia o al país

Finalmente se ignoran los que no corresponden a comarca, se registran aquellos en los que hay una fallo y con los que el registro corresponde a comarca se aplica otro *javascript* que convierte este registro en múltiples filas en la base de datos una por cada tipo y categoría de establecimiento y se normaliza el calculo de las plazas de camping según la regla del 2011

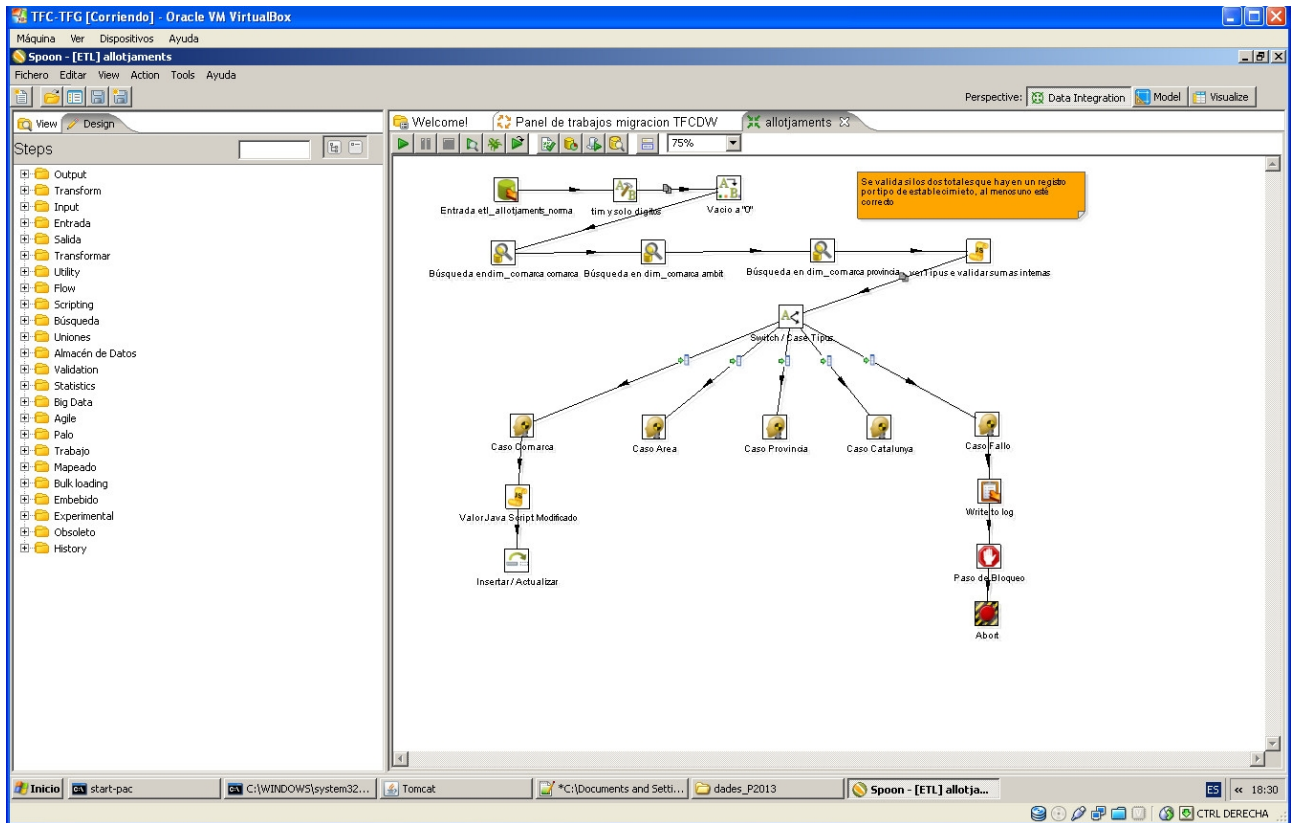


Ilustración 22: Transformación alojamientos descarte de filas.

#### 4.1.4. EQUIPAMIENTOS

En esta transformación con respecto al modelo lógico de migración se ha realizado cambios agrupando operaciones en busca de optimizar el tiempo de procesos ya hay que procesar 31771 registros.

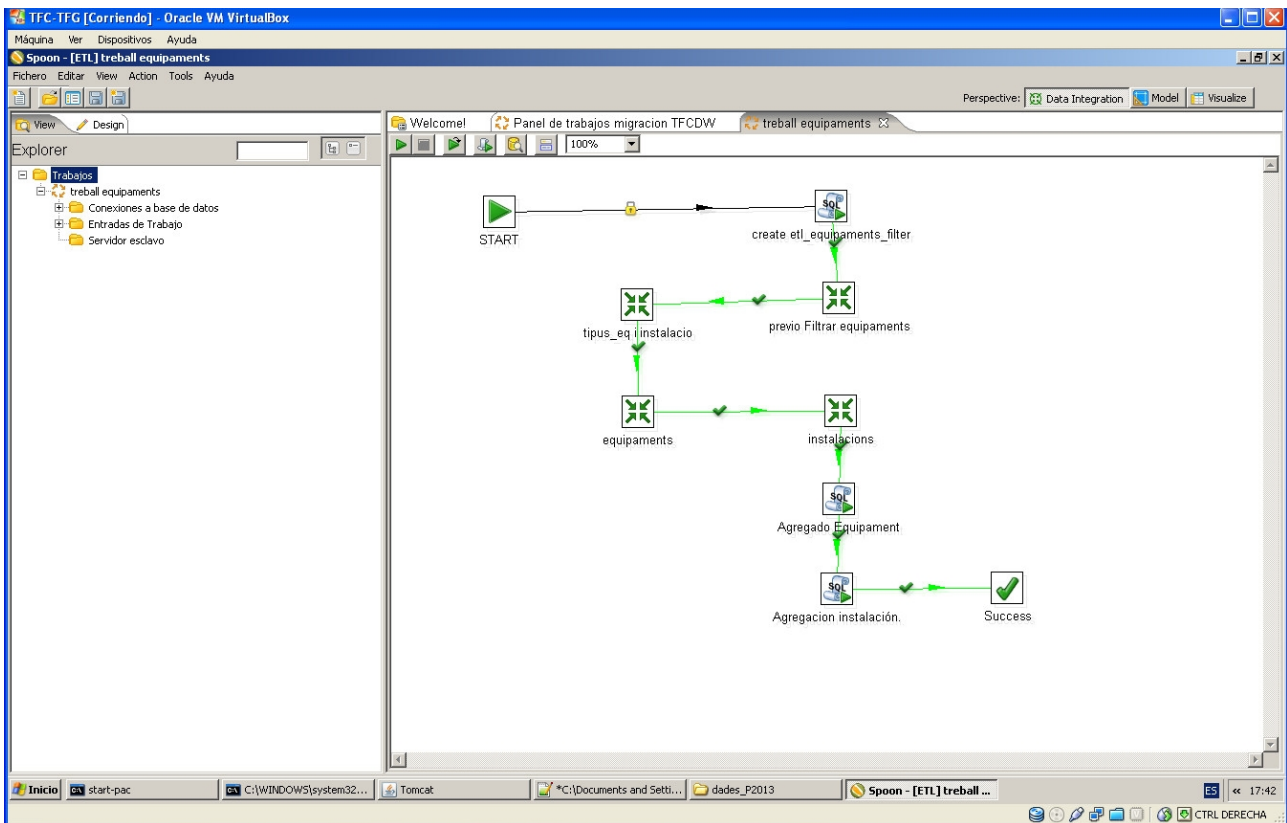


Ilustración 23: Transformación de equipamientos

Se realiza una transformación previa para retirar los registros duplicados y se le da al registro un valor único para ser referencia en todo el proceso de transformación, en ese instante también se les asigna las claves de municipio y comarca.

Durante dicha transformación se encontraron 149 de repetidos, para detectarlos se creó una verificación por *nom* y *cp* (código postal), más el resultado de aplicar md5 a categoría, para así poder tener en cuenta el caso de que un mismo equipamiento perteneciera a varios tipos de equipamiento.

El siguiente paso fue determinar los tipos de equipamiento e instalación, a partir de los campos categorías, del que se aplicó un *script* que recoge los 3 primeros campos de categoría como pertenecientes a la tabla de tipo de equipamiento (*familia*, *grup tipus*) y el restante para la tabla de instalación.

Una vez separados se aplicaron transformaciones para retirar las redundancias de familia y grupo y tipo por un lado y por otro las "instalaciones" para ser guardados en sus correspondientes dimensiones. A mayores se crean los ficheros intermediarios para que en posteriores pasos enlazar el equipamiento con los tipos de equipamiento e instalación que le corresponden.

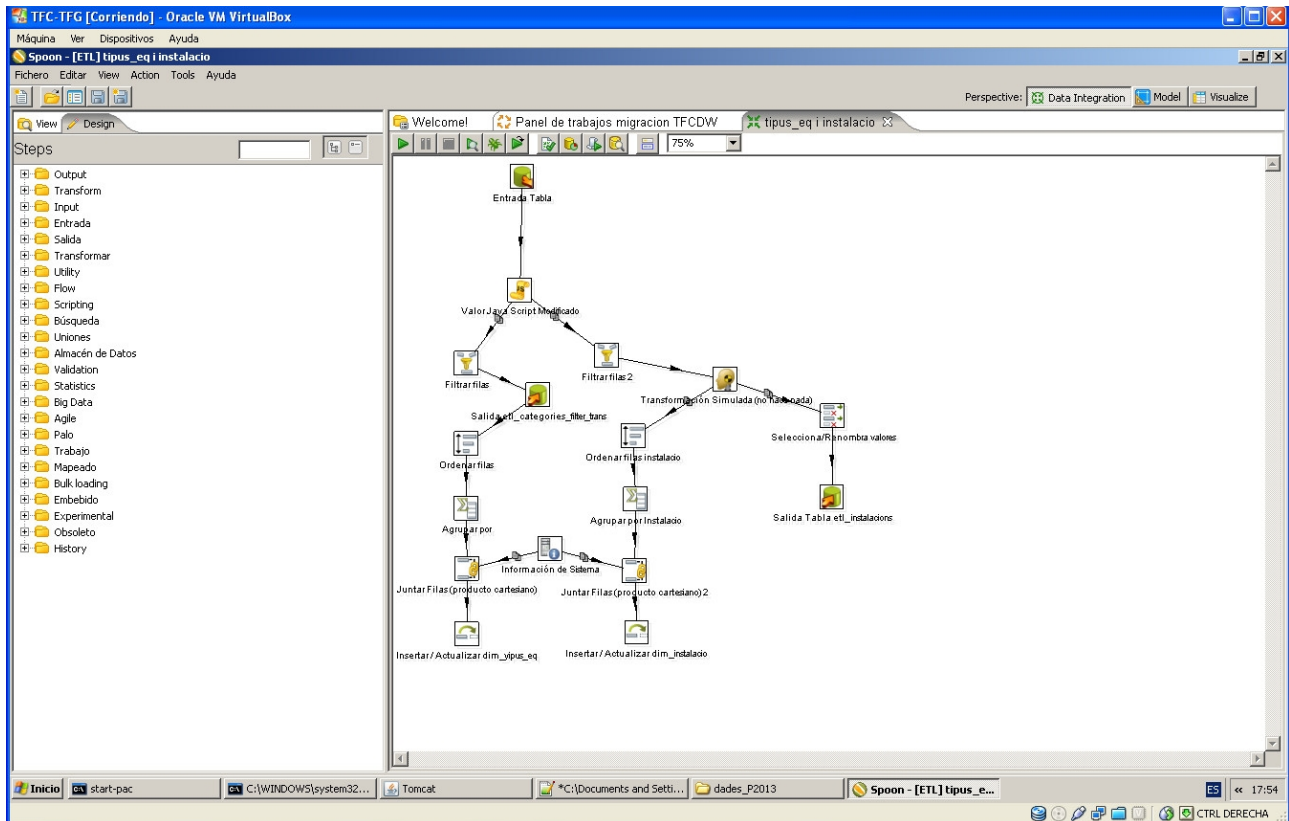


Ilustración 24: La generación de los categorías de equipamientos.

El siguiente paso *equipaments* es el encargado de guardar las tablas de *fct equipaments*, *deg equipaments* (en *deg equipaments* estará los datos de dirección del equipamiento)

Durante el proceso se detectaron 75 equipamientos que tenían mas de un registro con al misma combinación de tipo de equipamiento, eso quería decir que la diferencia estaba en las categorías que pertenecen a *instalació*. Por lo tanto han sido ignorados, para la carga en equipamientos, pero siguen el proceso de transformación para asignar dichas instalaciones al equipamiento en cuestión.

El siguiente paso es a partir de los ficheros intermedios generados en los pasos anteriores asignar a cada equipamiento relación con su lista de instalaciones.

Finalmente se lanzan los *scripts SQL* para que se calcule los ficheros de agregación de equipamientos e instalación con la suma de equipamientos por comarca.

## 4.2. PROCESO DE CREACIÓN DE LAS CONSULTAS OLAP Y DE LOS INFORMES.

El proceso se puede dividir en dos grandes secciones , aunque en realidad están intimamente relacionadas, por una lado el proceso de creación de los *metadatos* y por otro la creación de los informes.

### 4.2.1. PROCESO DE CREACIÓN DE LOS METADATOS Y ESQUEMAS.

Para el proceso de generación de los cubos de las consultas OLAP y de los informes, se intento

definir una única estructura única de *metadatos* y *esquemas OLAP*, común para todos los informes y consultas.

Los metadatos y los esquemas tienen por funcionalidad ser una capa intermedia entre los informes y la base de datos, que permite añadir información sobre los datos a mayores de la que se puede definir en el DDL de las Bases de datos.

El *esquema* es utilizado por el servidor **Mondrian ROLAP** y los *metadatos* son utilizados por el gestor de informes. Cada uno tienen sus propias gestiones de configuración y a su vez ambas estructuras se pueden crear desde la consola de **Pentaho BI** cuando se define una fuente de datos.

Cuales son las diferencias de crearlos en la consola del **Pentaho BI** o en los gestores específicos:

- Los gestores específicos el **Pentaho Metadata Editor** para los informes y el **Pentaho Schema Workbench** para los esquemas OLAP a diferencia de la consola permiten definir la estructura como un todo, es decir, permite definir elementos que después son reutilizados en las distintos informes. Así por ejemplo se define una sola vez la dimensión territorial.
- La ventaja de la consola del **Pentaho BI** es que ofrece una configuración simple, fácil y de uso muy intuitivo.

Pero la idea se abandonó ya que el esfuerzo para entrar los metadatos y esquemas, no compensaba con el uso posterior que se hacía de ellos, concretamente:

- Para la elaboración de los informes tanto los metadatos como los esquemas no valían de por sí y se habían precisado de utilizar la escritura de consultas en formatos **XML** o **MDX**, por lo tanto representaba el mismo esfuerzo que hacer una consulta **SQL**.
- Los informes y cubos requeridos son de información cruzada entre las diferentes estrellas por lo que la única manera de hacerlo es que se definan los esquemas a partir de consultas **SQL**.
- La ventaja de tener una sola definición de los datos utilizando las herramientas específicas en este almacén de datos no era representativa por el escaso número de jerarquías, por lo tanto no es problema duplicar la configuración varias veces.

Así que se decidió simplemente aplicar una convención de nombres a ser utilizado en cada configuración de fuente de datos en **Pentaho BI Server**.

- Dimensión "*Temporal*" para *dim\_any*
  - Jerarquía *Anual* {*anual*}
  - Jerarquía *Antiguetat* {*actual*}.
- Dimensión "*Territorial*" desde comarca.
  - Jerarquía territorial {*Provincia, Ambit, Comarca* }
- Dimensión "*Territorial*" desde municipio.
  - Jerarquía territorial {*Provincia, Ambit, Comarca, Municipio* }



- La dimensión tipos de equipamiento: “*Categoría*”
  - Jerarquía *categoría* {*Familia, Grup, Tipus* }
- Dimensión categorías de establecimiento: “*Categoría*”
  - Jerarquía *categoría* {*Tipus, Grup* }

#### 4.2.2. PROCESO DE CREACIÓN DE LOS INFORMES

Para realizar los listado primero se intentaron con la herramienta **Web Based Ad Hoc Query and Reporting** disponible en la consola de **Pentaho BI** pero nos encontramos con el problema desde las opciones de pantalla no permite incluir datos calculados, así que se descarto. En todo caso, se aprovecho para realiza una maquetación aceptable de los informes a modo de prototipo para después incluir los datos que faltan en el generador de informes **Pentaho Report Designer**

Con el **Pentaho Report Designer** se esperaba que con la importación de los dos prototipos al gestor de informes también incluyera la información de los metadatos y de la fuente de datos, pero eso no fue así, se probaron otras opciones como realizar conexiones vía metadatos y esquemas OLAP pero tampoco se ajustó a lo que se necesitaba, así que finalmente se utilizo como fuentes de datos las consultas SQL.

Para consulta al cubo OLAP se utilizo el frontal **JPivot** que accede a los cubos de **Mondrian Server** disponible en la consola **Pentaho BI**, existe un segundo frontal **Saiku**, pero no permite entrar campos calculados a partir del formato **MDX** y después, a mayores, se encontró con el problema que una vez guardado daba errores internos en el log de la consola al volverlo a cargar y olvidaba lo realizado

### 4.3. CONFIGURACIONES Y ACCESOS

Para simular un usuario final se ha creado el usuario *l'Observatori Nacional d'Ocupació ondo* (clave **ondo**) con el perfil de **ceo**

Se han asignado permisos de ejecución y poder programar procesos en diferido para los consultas OLAP creadas y para el caso de los ejemplos de acceso a los tres cubos (*estrellaDemograficTerritorial, estrellaEquipaments, estrellaEstabliments*) tiene todos los permisos. (de echo no he podido restringir los permisos de los cubos generados por **Saiku**)

El caso de los informes no hay ningún tipo de permisos al no poder publicarlos dentro del **Servidor BI**

Las bases de datos creadas son la *tfcetl* para la **Stagging Area** y la *tfcdw* para el almacén de datos.: `root > tfcmmd`

## 5. Presentación de los informes realizados.

---

### 5.1. CONSULTAS OLAP

Se han creado los siguientes cubos que se acceden desde la consola **Pentaho BI Server** y estan en la carpeta TFCDW, han sido creados por el usuario *tfcmmd* y usuario *ondo* puede consultarlos al tener un perfil de *ceo*.



*Observación: me he encontrado con muchos problemas, con las consultas OLAP ya que en ocasiones desaparecen o el servidor no responde. Y en ciertas ocasiones el usuario ondo no ve todas las consultas a las que tiene acceso. Supongo que son problemas debido a los recursos escasos de memoria y procesador.*

### 5.1.1. CONSULTA ESTABLIMENTS COMARQUES RATIOS

Cumple con los requerimientos de tener las siguientes ratios.

- Total d'establiments
- Total de places
- Oferta mitjana de places
- % de places respecte població
- Indicador de places vs persones
- Quantitat de places ofertes / superfície del territori
- Nombre d'establiments/Nombre d'equipaments

Se puede navegar en territorial, hasta nivel de comarca y por las categorías de los establecimientos.

Los datos de habitantes y equipamientos no son validos en todas las combinación del cubo ya que el sistema no permite elegir la función de agregación según la dimensión que se intervenga, es decir habitantes y equipamientos solo deberían de sumar a nivel territorial, a nivel de años debería de realizar promedio y con tipos de equipamiento no debería de hacer nada.

La fuente de datos utilizada es una consulta SQL.

The screenshot shows the Pentaho User Console interface with a pivot table titled 'EquipamentsComarcaRatios'. The table displays data for various regions (territorial) and categories (campings, hotels, trural) across the years 2005 and 2006. The columns represent different metrics: s\_places, s\_establiments, poblacio, equipament, Mitjana Places, Poblacio places, places habitants, equipaments establiment, and super. The data is organized in a grid format with rows for each region and category, and columns for each year and metric.

territorial	categories	s_places		s_establiments		poblacio		equipament		Mitjana Places		Poblacio places		places habitants		equipaments establiment		super
		2012		2006		2012		2006		2012		2006		2012		2006		2012
		2012	2006	2012	2006	2012	2006	2012	2006	2012	2006	2012	2006	2012	2006	2012	2006	2012
All territorials	campings	277,069	289,189	353	354	7,347,272	6,815,100	30,874	30,874	785	817	27	24	3.77%	4.24%	87.46	87.21	16,09
	hotels	294,777	259,120	2,837	2,624	7,347,272	6,815,100	30,874	30,874	104	99	25	26	4.01%	3.80%	10.88	11.77	16,09
	trural	16,968	12,865	2,156	1,669	7,347,272	6,815,100	30,874	30,874	8	8	433	530	0.23%	0.19%	14.32	18.50	16,09
Barcelona	campings	52,075	58,344	88	88	5,603,431	5,282,215	18,916	18,916	592	663	108	91	0.93%	1.10%	214.95	214.95	5,31
	hotels	138,460	110,756	1,234	1,083	5,603,431	5,282,215	18,916	18,916	112	102	40	48	2.47%	2.10%	15.33	17.47	5,31
	trural	5,339	3,859	674	512	5,603,431	5,282,215	18,916	18,916	8	8	1,050	1,369	0.10%	0.07%	28.07	36.95	5,31
Girona	campings	133,485	135,064	139	136	690,969	604,581	4,396	4,396	960	993	5	4	19.32%	22.34%	31.63	32.32	2,88
	hotels	79,100	77,507	818	805	690,969	604,581	4,396	4,396	97	96	9	8	11.45%	12.82%	5.37	5.46	2,88
	trural	5,426	3,777	659	465	690,969	604,581	4,396	4,396	8	8	127	160	0.79%	0.62%	6.67	9.45	2,88
Lleida	campings	22,224	22,498	60	61	381,831	340,131	3,629	3,629	370	369	17	15	5.82%	6.61%	60.48	59.49	4,16
	hotels	21,204	19,615	437	410	381,831	340,131	3,629	3,629	49	48	18	17	5.55%	5.77%	8.30	8.85	4,16
	trural	3,875	3,350	536	461	381,831	340,131	3,629	3,629	7	7	99	102	1.01%	0.98%	6.77	7.87	4,16
Tarragona	campings	69,285	73,283	66	69	671,041	588,173	3,933	3,933	1,050	1,062	10	8	10.33%	12.46%	59.59	57.00	3,74
	hotels	56,013	51,242	348	326	671,041	588,173	3,933	3,933	161	157	12	11	8.35%	8.71%	11.30	12.06	3,74
	trural	2,328	1,879	287	231	671,041	588,173	3,933	3,933	8	8	288	313	0.35%	0.32%	13.70	17.03	3,74

Il·lustració 25: Ratios Establecimientos comarca lado izquierdo

The screenshot shows the Pentaho User Console interface with a pivot table titled 'EquipamentsComarcaRatios'. The table displays data for various regions (territorial) and categories (campings, hotels, trural) across the years 2005 and 2006. The columns represent different metrics: s\_places, s\_establiments, poblacio, equipament, Mitjana Places, Poblacio places, places habitants, equipaments establiment, and superficie. The data is organized in a grid format with rows for each region and category, and columns for each year and metric. A warning message is displayed: 'JPivot has been replaced by Pentaho Analyzer. It is provided as a convenience but will no longer be enhanced or officially supported by Pentaho.'

territorial	categories	s_places		s_establiments		poblacio		equipament		Mitjana Places		Poblacio places		places habitants		equipaments establiment		superficie
		2012		2006		2012		2006		2012		2006		2012		2006		2012
		2012	2006	2012	2006	2012	2006	2012	2006	2012	2006	2012	2006	2012	2006	2012	2006	2012
All territorials	campings	277,069	289,189	353	354	7,347,272	6,815,100	30,874	30,874	785	817	27	24	3.77%	4.24%	87.46	87.21	16,093
	hotels	294,777	259,120	2,837	2,624	7,347,272	6,815,100	30,874	30,874	104	99	25	26	4.01%	3.80%	10.88	11.77	16,093
	trural	16,968	12,865	2,156	1,669	7,347,272	6,815,100	30,874	30,874	8	8	433	530	0.23%	0.19%	14.32	18.50	16,093

Il·lustració 26: Ratios Establecimientos comarca lado derecho

### 5.1.2. CONSULTA ESTABLIMENTS GÈNERE TURISME

Cumple con los requerimientos de tener los siguientes ratio.

- Total d'establiments
- Total de places
- Indicador d'establiments vs habitants per gènere
- Nombre d'establiments/Nombre d'equipaments

Ademas muestra el total por género, el total de la población y en vez de mostrar todos los equipamientos se muestra una selección de ellos bajo el criterio que tengan que ver con turismo: *Esport i lleura, cultura, turisme*. Y finalmente un ratio el % de mujeres.

Observación en teoría el ratio del % de mujeres debería de ser constante en todos los años, pero hay diferencias decimales por el redondeo que se ha realizado para calcularlos.

Los datos de habitantes y equipamientos no son válidos en todas las combinaciones del cubo ya que el sistema no permite elegir la función de agregación según la dimensión que se intervenga, es decir habitantes y equipamientos solo deberían de sumar a nivel territorial, a nivel de años debería de realizar promedio y con tipos de equipamiento no debería de hacer nada.

La fuente de datos utilizada es una consulta SQL.

territorial	tipus	poblacio		dones		homes		s_places		s_establiments		esport_i_lleura		cultura		turisme	
		2012	2006	2012	2006	2012	2006	2012	2006	2012	2006	2012	2006	2012	2006	2012	2006
All territorials	campings	7,347,272	6,815,100	3,631,869	3,364,785	3,715,403	3,450,315	277,069	289,189	353	354	10,057	10,057	2,699	2,699	219	219
	hotels	7,347,272	6,815,100	3,631,869	3,364,785	3,715,403	3,450,315	294,777	259,120	2,837	2,624	10,057	10,057	2,699	2,699	219	219
	trural	7,347,272	6,815,100	3,631,869	3,364,785	3,715,403	3,450,315	16,968	12,865	2,156	1,669	10,057	10,057	2,699	2,699	219	219
Barcelona	campings	5,603,431	5,282,215	2,752,850	2,592,705	2,850,581	2,689,510	52,075	58,344	88	88	5,894	5,894	1,565	1,565	68	68
	hotels	5,603,431	5,282,215	2,752,850	2,592,705	2,850,581	2,689,510	138,460	110,756	1,234	1,083	5,894	5,894	1,565	1,565	68	68
	trural	5,603,431	5,282,215	2,752,850	2,592,705	2,850,581	2,689,510	5,339	3,859	674	512	5,894	5,894	1,565	1,565	68	68
Girona	campings	690,969	604,581	347,567	303,905	343,402	300,676	133,485	135,064	139	136	1,567	1,567	443	443	55	55
	hotels	690,969	604,581	347,567	303,905	343,402	300,676	79,100	77,507	818	805	1,567	1,567	443	443	55	55
	trural	690,969	604,581	347,567	303,905	343,402	300,676	5,426	3,777	659	465	1,567	1,567	443	443	55	55
Lleida	campings	381,831	340,131	193,428	172,268	188,403	167,863	22,224	22,498	60	61	1,290	1,290	311	311	45	45
	hotels	381,831	340,131	193,428	172,268	188,403	167,863	21,204	19,615	437	410	1,290	1,290	311	311	45	45
	trural	381,831	340,131	193,428	172,268	188,403	167,863	3,875	3,350	536	461	1,290	1,290	311	311	45	45
Tarragona	campings	671,041	588,173	338,024	295,907	333,017	292,266	69,285	73,283	66	69	1,306	1,306	380	380	51	51
	hotels	671,041	588,173	338,024	295,907	333,017	292,266	56,013	51,242	348	326	1,306	1,306	380	380	51	51
	trural	671,041	588,173	338,024	295,907	333,017	292,266	2,328	1,879	287	231	1,306	1,306	380	380	51	51

Ilustración 27: Volcado de la pantalla de la consulta establecimientos género lado izquierdo

The screenshot shows the Pentaho User Console interface. The main window displays a data table with the following columns: s\_places, s\_establiments, esport\_lleure, cultura, turisme, Establiments dones, Establiments homes, dones poblacio, and places equipament. Each column has sub-columns for the years 2006 and 2012. The table is rotated 90 degrees clockwise. Below the table, a query log shows a SQL query: SELECT \* FROM fctdw\_fct\_demografid LIMIT 0, 1000. The status indicates 1000 row(s) returned.

	s_places		s_establiments		esport_lleure		cultura		turisme		Establiments dones		Establiments homes		dones poblacio		places equipament		
	2006	2012	2006	2012	2006	2012	2006	2012	2006	2012	2006	2012	2006	2012	2006	2012	2006	2012	
3	3,450,315	277,069	289,189	353	354	10,057	10,057	2,699	2,699	219	219	0.01%	0.01%	0.01%	0.01%	49.43%	49.37%	21.35	22.29
8	3,450,315	294,777	259,120	2,837	2,624	10,057	10,057	2,699	2,699	219	219	0.08%	0.08%	0.08%	0.08%	49.43%	49.37%	22.72	19.97
8	3,450,315	16,968	12,865	2,156	1,669	10,057	10,057	2,699	2,699	219	219	0.06%	0.05%	0.06%	0.05%	49.43%	49.37%	1.31	0.99
1	2,689,510	52,075	58,344	88	88	5,894	5,894	1,565	1,565	68	68	0.00%	0.00%	0.00%	0.00%	49.13%	49.08%	6.92	7.75
1	2,689,510	138,460	110,756	1,234	1,083	5,894	5,894	1,565	1,565	68	68	0.04%	0.04%	0.04%	0.04%	49.13%	49.08%	18.40	14.71
1	2,689,510	5,399	3,859	674	512	5,894	5,894	1,565	1,565	68	68	0.02%	0.02%	0.02%	0.02%	49.13%	49.08%	0.71	0.51
2	300,676	133,485	135,064	139	136	1,567	1,567	443	443	55	55	0.04%	0.04%	0.04%	0.05%	50.30%	50.27%	64.64	65.41
2	300,676	79,100	77,507	818	805	1,567	1,567	443	443	55	55	0.24%	0.26%	0.24%	0.27%	50.30%	50.27%	38.31	37.53
2	300,676	5,426	3,777	659	465	1,567	1,567	443	443	55	55	0.19%	0.15%	0.19%	0.15%	50.30%	50.27%	2.63	1.83
8	167,863	22,224	22,498	60	61	1,290	1,290	311	311	45	45	0.03%	0.04%	0.03%	0.04%	50.66%	50.65%	13.50	13.67
8	167,863	21,204	19,615	437	410	1,290	1,290	311	311	45	45	0.23%	0.24%	0.23%	0.24%	50.66%	50.65%	12.88	11.92
8	167,863	3,875	3,350	536	461	1,290	1,290	311	311	45	45	0.28%	0.27%	0.28%	0.27%	50.66%	50.65%	2.35	2.04
7	292,266	69,285	73,283	66	69	1,306	1,306	380	380	51	51	0.02%	0.02%	0.02%	0.02%	50.37%	50.31%	39.89	42.19
7	292,266	56,013	51,242	348	326	1,306	1,306	380	380	51	51	0.10%	0.11%	0.10%	0.11%	50.37%	50.31%	32.25	29.50
7	292,266	2,328	1,879	287	231	1,306	1,306	380	380	51	51	0.08%	0.08%	0.09%	0.08%	50.37%	50.31%	1.34	1.08

Il·lustració 28: Volcado de la pantalla de la consulta establecimientos género lado derecho



### 5.1.3. CONSULTA EQUIPAMENTS COMARCA RATIOS

Cumple con los requerimientos de tener los siguientes ratio.

- Nombre d'establiments/Nombre d'equipaments
- % de població per equipament
- Indicador d'equipaments vs població

Ademas muestra los agregados de equipamientos, habitantes y establecimientos y se puede navegar por las dimensiones anual, tipos de equipamientos y territorial hasta comarca.

Los datos de habitantes y establecimientos no son validos en todas las combinación del cubo ya que el sistema no permite elegir la función de agregación según la dimensión que se intervenga, es decir habitantes y establecimientos solo deberían de sumar a nivel territorial, a nivel de años debería de realizar promedio y con tipos de equipamiento no debería de hacer nada.

La fuente de datos utilizada es una consulta SQL.

	Measures	equipaments		habitants		establiments		Establiments equipament		habitants/ equipament		equipament / habitants	
		anual	anual	2012	2006	2012	2006	2012	2006	2012	2006	2012	2006
tipus equipaments	Territorial	2012	2006	2012	2006	2012	2006	2012	2006	2012	2006	2012	2006
Administracio_Publica	All Territorials	828	828	7,347,272	6,815,100	5,346	4,647	6.457	5.612	8873.5	8230.8	0.011%	0.012%
Agricultura_Ramaderia_Pesca	All Territorials	88	88	7,347,272	6,815,100	5,346	4,647	60.750	52.807	83491.7	77444.3	0.001%	0.001%
Associacionisme_participacio	All Territorials	13	13	7,347,272	6,815,100	5,346	4,647	411.231	357.462	565174.8	524238.5	0.000%	0.000%
Comerc_Consum	All Territorials	351	351	7,347,272	6,815,100	5,346	4,647	15.231	13.239	20932.4	19416.2	0.005%	0.005%
Cultura	All Territorials	2,699	2,699	7,347,272	6,815,100	5,346	4,647	1.981	1.722	2722.2	2525.0	0.037%	0.040%
Economia	All Territorials	50	50	7,347,272	6,815,100	5,346	4,647	106.920	92.940	146945.4	136302.0	0.001%	0.001%
Educacio	All Territorials	5,347	5,347	7,347,272	6,815,100	5,346	4,647	1.000	0.869	1374.1	1274.6	0.073%	0.078%
Emergencies_seguretat	All Territorials	433	433	7,347,272	6,815,100	5,346	4,647	12.346	10.732	16968.3	15739.3	0.006%	0.006%
Empresa_Industria_energia	All Territorials	61	61	7,347,272	6,815,100	5,346	4,647	87.639	76.180	120447.1	111723.0	0.001%	0.001%
Esport_Lleure	All Territorials	10,057	10,057	7,347,272	6,815,100	5,346	4,647	0.532	0.462	730.6	677.6	0.137%	0.148%
Habitatge	All Territorials	96	96	7,347,272	6,815,100	5,346	4,647	55.688	48.406	76534.1	70990.6	0.001%	0.001%
Justicia	All Territorials	1,602	1,602	7,347,272	6,815,100	5,346	4,647	3.337	2.901	4586.3	4254.1	0.022%	0.024%
Llengua_Comunicacio	All Territorials	376	376	7,347,272	6,815,100	5,346	4,647	14.218	12.359	19540.6	18125.3	0.005%	0.006%
Medi_ambient	All Territorials	376	376	7,347,272	6,815,100	5,346	4,647	14.218	12.359	19540.6	18125.3	0.005%	0.006%
Mobilitat_Transports	All Territorials	186	186	7,347,272	6,815,100	5,346	4,647	18.158	15.681	23111.7	21088.1	0.003%	0.003%

Ilustración 29: Volcado de pantalla con la consulta equipamientos comarca

### 5.1.4. CONSULTAS OLAP DE LOS DATA-MARTS.

A mayores se ha creado las fuentes de conexión y los esquemas de análisis para poder acceder mediante consultas OLAP a las tres estrellas definidas, el usuario tiene libertad de crear las consultas que deseé, se han dejado accesibles 3 consultas. Han sido creadas con Saiku por que es una herramienta intuitiva y apta para el usuario, aunque, tiene un problema ya que ciertas consultas guardadas en Saiku después no son recuperables, sobre todo si se realizan con varios niveles de

profundidad.

Las consultas son:

- estrellaDemograficTerritorial Que hace uso de la fuente/esquema DemograficTerritorial, que permite navegar hasta nivel de municipio.
- estrellaEquipaments Que hace uso de la fuente/esquema: estrellaEquipaments, que permite navegar hasta nivel de municipio
- estrellaEstabliments, Que hace uso de la fuente/esquema: estrellaEstabliments

En las capturas de pantalla a parte de la consulta se puede apreciar la estructuras de las dimensiones y medidas.

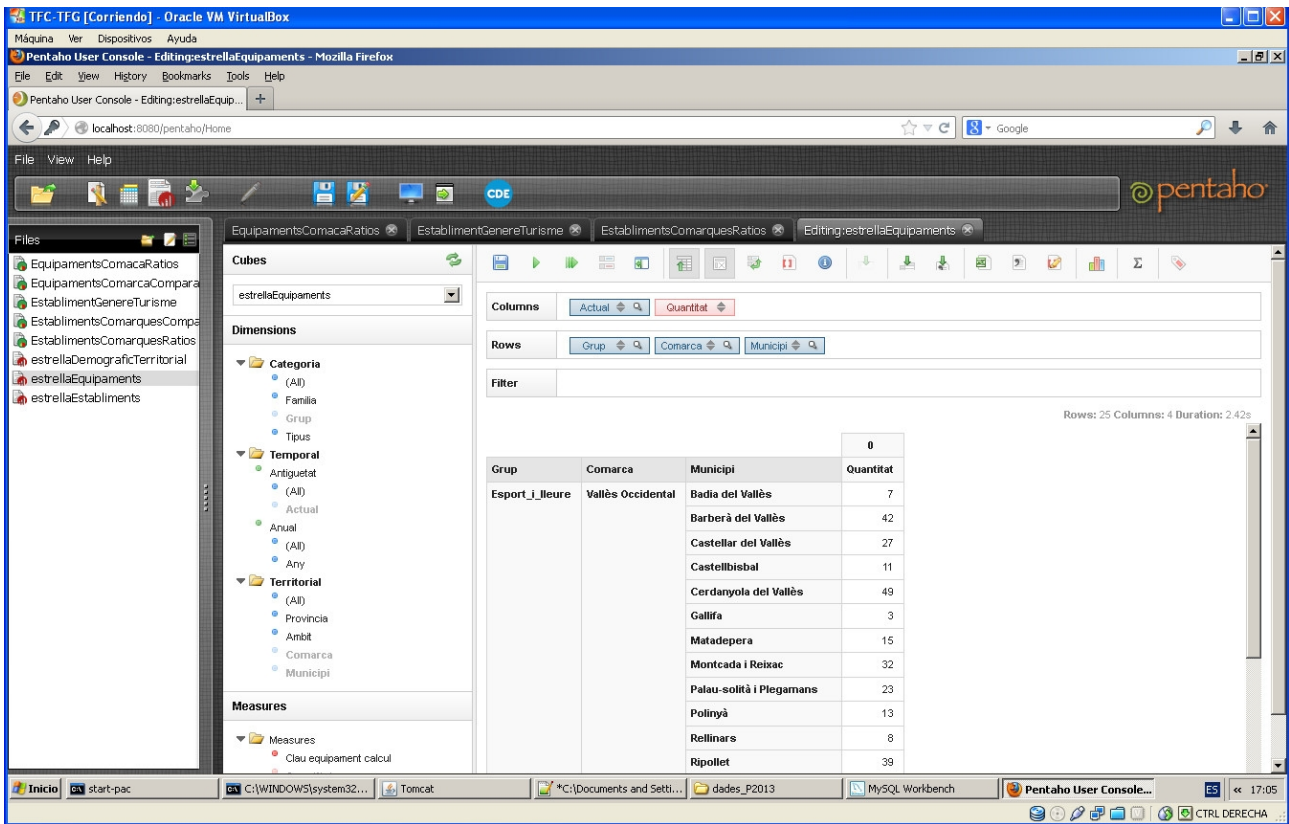
The screenshot shows the Pentaho User Console interface. On the left, a tree view displays the cube structure with dimensions: Temporal (Anual, Clau any, Antiguetat, Actual) and Territorial (Provincia, Ambit, Comarca, Municipi). The main area shows a table with columns for 'Extension' and 'Habitants' across years 2006-2012. The table data is as follows:

Provincia	Extension						Habitants						
	2006	2007	2008	2009	2010	2011	2006	2007	2008	2009	2010	2011	
Barcelona	5.310	5.310	5.310	5.310	5.310	5.310	5.282.215	5.370.355	5.397.801	5.485.867	5.580.255	5.584.138	5.6
Girona	2.880	2.880	2.880	2.880	2.880	2.880	604.581	626.839	644.093	667.712	682.770	687.798	6
Lleida	4.160	4.160	4.160	4.160	4.160	4.160	340.131	347.796	353.925	366.672	375.784	379.147	3
Tarragona	3.736	3.736	3.736	3.736	3.736	3.736	588.173	607.406	628.544	654.243	685.139	669.367	6

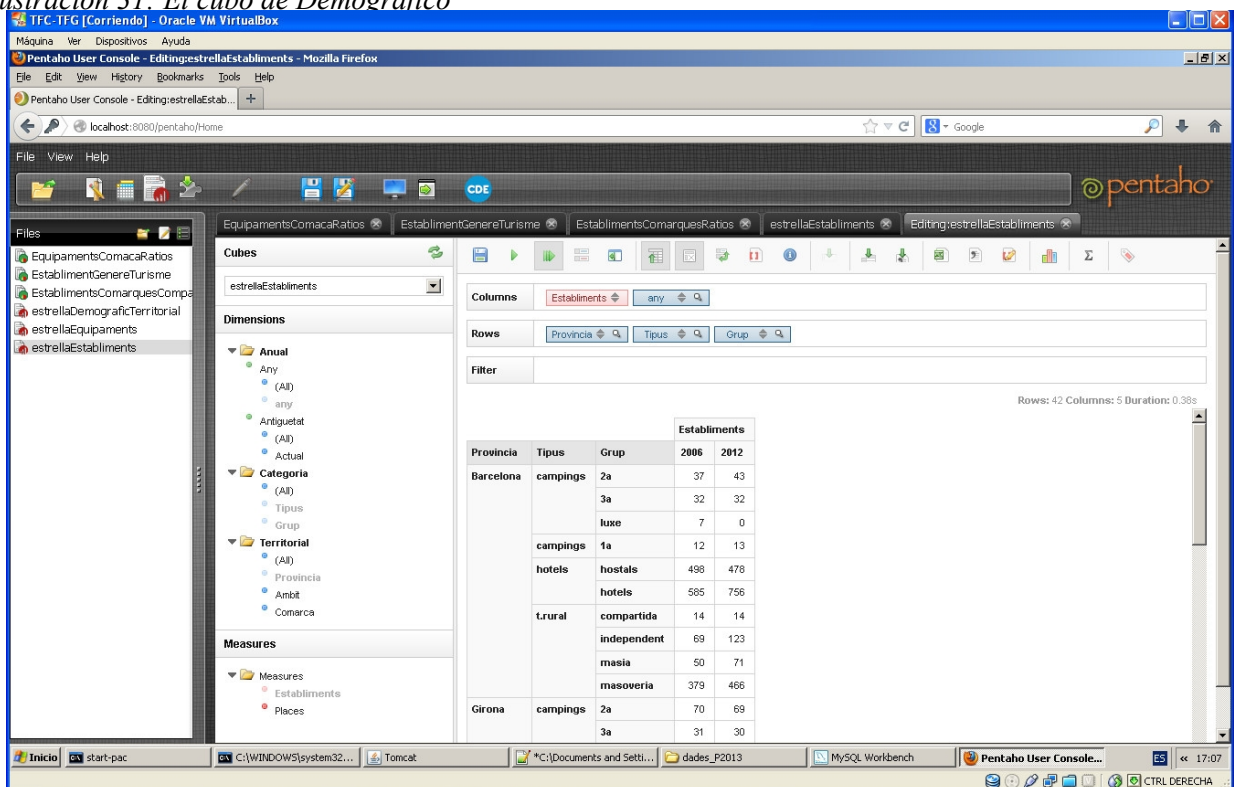
The bottom of the screenshot shows a query log with the following entry:

```
12 16:47:53 SELECT * FROM ftcw_ftc_demografic LIMIT 0, 1000 1000 row(s) returned 0.010 sec / 0.000 sec
```

Ilustración 30: El cubo de equipamientos.



Il·lustració 31: El cubo de Demogràfico



Il·lustració 32: El cubo de los establecimientos.

## **5.2. INFORMES**

Se han creado los siguientes informes, que se encuentran en el directorio *mis documentos/* a continuación se detalla que datos de los solicitados en el proyecto contiene cada uno de ellos.

Observación: no se ha podido incluirlos en la consola del **Pentaho BI Server**

### **5.2.1. INFORME ESTABLIMENTS COMARCA (INFORMEESTABLIMENTCOMARCA.PRPT)**

Cumple con los requerimientos de tener los siguientes ratio.

- Total d'establiments
- Total de places
- Oferta mitjana de places
- % de places respecte població
- Indicador de places vs persones
- Quantitat de places ofertes / superfície del territori

Con detalle a nivel de comarca y acumulados a nivel de provincia, tipo de establecimiento (camping,hotel,t. rural) y año.



Informe Establiments Comarca

mayo 29, 2013 @ 05:14

Anys  
Anys: 2006

Tipus: campings  
 Provincia: Barcelona

Ambit	Comarca	Establiments	Places	Mitjana	Poblacio	% places persones per poblacio	Superficie Territori	Places / Km2
Comarques Centrals	Bages	3	504	168	164.085	0,3%	325,57	834 0,60
Comarques Centrals	Berguedà	14	5.241	374	34.564	15,2%	6,59	335 15,64
Comarques Centrals	Osona	7	2.343	335	133.673	1,6%	57,05	548 4,28
Comarques Centrals	Solsonès	5	2.308	462	9.540	24,2%	4,13	22 104,91
metropoli de Barcelona	Baix Llobregat	3	8.339	2.780	757.814	1,1%	90,88	486 17,16
metropoli de Barcelona	Barcelonès	0	0	0	2.215.581	0%	,00	144 0,00
metropoli de Barcelona	Maresme	33	21.461	650	398.187	5,4%	18,55	394 54,47
metropoli de Barcelona	Valles Occidental	0	0	0	814.813	0%	,00	546 0,00
metropoli de Barcelona	Valles Oriental	10	4.268	427	357.500	1,2%	83,76	642 6,65
Penedès	Alt Penedès	1	376	376	89.902	0,4%	238,10	549 0,68
Penedès	Anoia	1	76	76	99.616	0,1%	1.310,76	361 0,21
Penedès	Baix Penedès	5	4.856	971	79.910	6,1%	16,27	264 19,39
Penedès	Garraf	6	8.572	1.429	127.926	6,7%	14,92	185 46,34
<b>Total Barcelona (campings)</b>		<b>88</b>	<b>58.344</b>	<b>663</b>	<b>5.282.215</b>	<b>1,1%</b>	<b>90,54</b>	<b>5.310 10,99</b>

Provincia: Girona

Ambit	Comarca	Establiments	Places	Mitjana	Poblacio	% places persones per poblacio	Superficie Territori	Places / Km2
Comarques gironines	Alt Empordà	33	38.621	1.170	96.699	39,1%	2,56	447 96,40
Comarques gironines	Baix Empordà	44	62.008	1.409	113.926	54,4%	1,84	509 121,82
Comarques gironines	Garrotxa	16	3.482	218	47.096	7,4%	13,53	382 9,12
Comarques gironines	Gironès	2	1.216	608	155.251	0,8%	127,67	364 3,34
Comarques gironines	La Selva	25	23.593	944	142.152	16,6%	6,03	767 30,76
Comarques gironines	Pla de l'Estany	4	2.017	504	25.150	8%	12,47	107 16,85
Comarques gironines	Ripollès	12	4.127	344	22.305	18,5%	5,40	304 13,88
<b>Total Girona (campings)</b>		<b>136</b>	<b>135.064</b>	<b>993</b>	<b>604.581</b>	<b>22,3%</b>	<b>4,48</b>	<b>2.880 46,90</b>

Provincia: Lleida

Ambit	Comarca	Establiments	Places	Mitjana	Poblacio	% places persones per poblacio	Superficie Territori	Places / Km2
Alt Llobregat	Alt Llobregat	0	3.606	776	16.106	46,8%	6,03	224 41,19

Il·lustració 34: El detall del informe establecimientos comarca

Informe Establiments Comarca

mayo 29, 2013 @ 05:14

Ponent	Segarra	46	309	7	18.396	1,7%	58,53	191 1,62
Ponent	Segrià	8	57	7	198.286	0%	3.478,70	955 0,06
Ponent	Urgell	35	287	8	32.579	0,9%	113,52	289 0,99
<b>Total Lleida (L.rural)</b>		<b>536</b>	<b>3.875</b>	<b>7</b>	<b>381.831</b>	<b>1%</b>	<b>98,54</b>	<b>4.160 0,93</b>

Provincia: Tarragona

Ambit	Comarca	Establiments	Places	Mitjana	Poblacio	% places persones per poblacio	Superficie Territori	Places / Km2
Camp de Tarragona	Alt Camp	49	412	8	37.645	1,1%	91,37	223 1,85
Camp de Tarragona	Baix Camp	34	275	8	186.271	0,1%	677,35	450 0,61
Camp de Tarragona	Costa de Barberà	39	296	8	17.662	1,3%	577,71	300 0,98
Camp de Tarragona	Priorat	48	415	9	3.978	10,4%	9,59	95 4,37
Camp de Tarragona	Tarragonès	10	73	7	248.529	0%	3.404,51	273 0,27
Terres de l'Ebre	Baix Ebre	41	263	6	78.938	0,3%	300,14	755 0,35
Terres de l'Ebre	Montsià	30	301	10	69.494	0,4%	230,88	581 0,52
Terres de l'Ebre	Ribera d'Ebre	14	102	7	19.016	0,5%	196,43	582 0,18
Terres de l'Ebre	Terra Alta	22	191	9	10.088	1,9%	52,82	484 0,38
<b>Total Tarragona (L.rural)</b>		<b>287</b>	<b>2.328</b>	<b>8</b>	<b>671.041</b>	<b>0,3%</b>	<b>288,25</b>	<b>3.743 0,62</b>
<b>Total L.rural</b>		<b>2.156</b>	<b>16.968</b>	<b>8</b>	<b>7.347.272</b>	<b>0,2%</b>	<b>433,01</b>	<b>16.093 1,05</b>
<b>Total 2012</b>		<b>5.346</b>	<b>588.814</b>	<b>110</b>	<b>7.347.272</b>	<b>8%</b>	<b>12,48</b>	<b>16.093 36,59</b>

Il·lustració 33: Los totales a nivel año tipo y provincia establecimientos comarca

### 5.2.2. INFORME ESTABLIMENTS COMARCA HOMES DONES EQUIPAMENTS (INFORMEESTABLIMENTSCOMARCA.B.PRPT)

Cumple con los requerimientos de tener los siguientes ratio.

- Total d'establiments
- Total de places
- Indicador d'establiments vs habitants per gènere
- Nombre d'establiments/Nombre d'equipaments

Con detalle a nivel de comarca y acumulados a nivel de provincia, tipo de establecimiento (camping,hotel,t. rural) y año.

mayo 29, 2013 @ 05:18

Anys  
Anys: 2006

Tipus: campings  
Provincia: Barcelona

Ambit	Comarca	Establiments	Places	Homes	Dones	homes / establiment	dones / establiment	Equipaments	equipaments / establiment
Comarques Centrals	Bages	3	504	82.656	81.429	27.552	27.143	1.066	355,3
Comarques Centrals	Berguedà	14	5.241	17.359	17.205	1.240	1.229	408	29,1
Comarques Centrals	Osona	7	2.343	66.972	66.701	9.567	9.529	996	142,3
Comarques Centrals	Solsonès	5	2.308	4.748	4.792	950	958	154	30,8
metropolità de Barcelona	Baix Llobregat	3	8.339	381.192	376.622	127.064	125.541	2.382	794,0
metropolità de Barcelona	Barcelonès	0	0	1.149.696	1.065.885	0	0	5.943	0,0
metropolità de Barcelona	Maresme	33	21.461	200.255	197.932	6.068	5.998	1.643	49,8
metropolità de Barcelona	Valles Occidental	0	0	411.413	403.400	0	0	2.619	0,0
metropolità de Barcelona	Valles Oriental	10	4.268	177.905	179.595	17.791	17.960	1.524	152,4
Penedès	Alt Penedès	1	376	44.674	45.228	44.674	45.228	662	662,0
Penedès	Anoia	1	76	49.503	50.115	49.503	50.115	621	621,0
Penedès	Baix Penedès	5	4.856	38.960	40.050	7.792	8.010	430	86,0
Penedès	Garraf	6	8.572	64.177	63.751	10.696	10.625	468	78,0
<b>Total Barcelona (campings)</b>		<b>88</b>	<b>58.344</b>	<b>2.689.510</b>	<b>2.592.705</b>	<b>30.563</b>	<b>29.463</b>	<b>18.916</b>	<b>215,0</b>
Provincia: Girona									
Ambit	Comarca	Establiments	Places	Homes	Dones	homes / establiment	dones / establiment	Equipaments	equipaments / establiment
Comarques gironines	Alt Empordà	33	38.621	49.012	49.687	1.485	1.506	949	28,8
Comarques gironines	Baix Empordà	44	62.008	56.900	57.028	1.293	1.296	779	17,7
Comarques gironines	Garrotxa	16	3.482	23.633	23.463	1.477	1.466	407	25,4
Comarques gironines	Gironès	2	1.216	77.974	77.277	38.987	38.639	984	492,0
Comarques gironines	La Selva	25	23.593	69.308	72.844	2.772	2.914	726	29,0
Comarques gironines	Baix Empordà	4	2.017	12.686	12.686	3.141	3.148	108	40,8

Ilustración 35: El detalle del informe Establiments Comarca Homes Dones Equipaments

Ambit	Comarca	Establiments	Places	Homes	Dones	homes / establiment	dones / establiment	Equipaments	equipaments / establiment
Provincia: Tarragona									
<b>Total Lleida ( t.rural )</b>		<b>536</b>	<b>3.875</b>	<b>188.403</b>	<b>193.428</b>	<b>351</b>	<b>361</b>	<b>3.629</b>	<b>6,8</b>
Provincia: Tarragona									
Camp de Tarragona	Alt Camp	49	412	18.496	19.149	377	391	339	6,9
Camp de Tarragona	Baix Camp	34	275	93.675	92.596	2.755	2.723	844	24,8
Camp de Tarragona	Conca de Barberà	39	296	8.433	8.649	216	222	258	6,6
Camp de Tarragona	Priorat	48	415	1.952	2.026	41	42	237	4,9
Camp de Tarragona	Tarragonès	10	73	123.376	125.153	12.338	12.515	965	96,5
Terres de l'Ebre	Baix Ebre	41	263	38.656	40.282	943	982	445	10,9
Terres de l'Ebre	Montsià	30	301	34.205	35.289	1.140	1.176	383	12,8
Terres de l'Ebre	Ribera d'Ebre	14	102	9.354	9.662	668	690	263	18,8
Terres de l'Ebre	Terra Alta	22	191	4.870	5.218	221	237	199	9,0
<b>Total Tarragona ( t.rural )</b>		<b>287</b>	<b>2.328</b>	<b>333.017</b>	<b>338.024</b>	<b>1.160</b>	<b>1.178</b>	<b>3.933</b>	<b>13,7</b>
<b>Total t.rural</b>		<b>2.156</b>	<b>16.968</b>	<b>3.715.403</b>	<b>3.631.869</b>	<b>1.723</b>	<b>1.723</b>	<b>30.874</b>	<b>14,3</b>
<b>Total 2012</b>		<b>5.346</b>	<b>588.814</b>	<b>3.715.403</b>	<b>3.631.869</b>	<b>695</b>	<b>679</b>	<b>30.874</b>	<b>5,8</b>

Il·lustració 36: Totals provincia, categoria i any informe Establiments Comarca Homes Dones Equipaments.

### 5.2.3. INFORME EQUIPAMENTS COMARCA (INFORMEEQUIPAMENTSCOMARCA.PRPT)

Cumple con los requerimientos de tener los siguientes ratio.

- Nombre d'establiments/Nombre d'equipaments
- % de població per equipament
- Indicador d'equipaments vs població

Con detalle a nivel de comarca y acumulados a nivel de provincia, a nivel de grupo de equipamiento (el nivel 2) ya que todos pertenecen a la misma familia de "equipamientos", y año

Muestra los equipamientos y su relación con habitantes y establecimientos (totalizados por comarca) y los ratios "equipaments/ poblacio %", "habitants / equipaments" y "establiment / equipament".



mayo 29, 2013 @ 05:24  
any: 2.006

grup: Administració\_Pública  
provincia: Barcelona

ambit	comarca	equipaments	habitants	establiments	equipaments / població	habitants / equipament	establiment / equipament
Comarques Centrals	Bages	33	164.085	96	0,02%	4.972,3	2,9
Comarques Centrals	Berguedà	16	34.564	160	0,05%	2.160,3	10,0
Comarques Centrals	Osona	44	133.673	168	0,03%	3.038,0	3,8
Comarques Centrals	Solsonès	6	9.540	112	0,06%	1.590,0	18,7
metropolità de Barcelo	..Baix Llobregat	27	757.814	77	0%	28.067,2	2,9
metropolità de Barcelo	..Barcelonès	9	2.215.581	464	0%	246.175,7	51,6
metropolità de Barcelo	..Maresme	27	398.187	228	0,01%	14.747,7	8,4
metropolità de Barcelo	..Vallès Occidental	21	814.813	38	0%	38.800,6	1,8
metropolità de Barcelo	..Vallès Oriental	40	357.500	87	0,01%	8.937,5	2,2
Penedès	Alt Penedès	27	89.902	61	0,03%	3.329,7	2,3
Penedès	Anoia	23	99.618	62	0,02%	4.331,2	2,7
Penedès	Baix Penedès	14	79.010	56	0,02%	5.643,6	4,0
Penedès	Garraf	6	127.928	74	0%	21.321,3	12,3
<b>Total Barcelona (Administració ...</b>		<b>293</b>	<b>5.282.215</b>	<b>1.683</b>	<b>0,01%</b>	<b>18.028,0</b>	<b>5,7</b>
provincia: Girona							
Comarques gironines	Alt Empordà	62	98.699	346	0,06%	1.591,9	5,6
Comarques gironines	Baix Empordà	32	113.928	277	0,03%	3.560,3	8,7
Comarques gironines	Garroba	19	47.096	150	0,04%	2.478,7	7,9
Comarques gironines	Gironès	25	155.251	63	0,02%	6.210,0	2,5
Comarques gironines	La Selva	23	142.152	335	0,02%	6.180,5	14,6
Comarques gironines	Pla de l'Estany	8	25.150	84	0,03%	3.143,8	10,5
Comarques gironines	Ripollès	19	22.305	151	0,09%	1.173,9	7,9
<b>Total Girona (Administració_Pub ...</b>		<b>188</b>	<b>604.581</b>	<b>1.406</b>	<b>0,03%</b>	<b>3.215,9</b>	<b>7,5</b>

Il·lustració 37: Cabecera y cuerpo del informe Equipaments Comarca

mayo 29, 2013 @ 05:24  
any: 2.008

grup: Administració\_Pública  
provincia: Barcelona

ambit	comarca	equipaments	habitants	establiments	equipaments / població	habitants / equipament	establiment / equipament
Comarques Centrals	Bages	33	171.654	103	0,02%	5.201,6	3,1
Comarques Centrals	Berguedà	16	35.234	181	0,05%	2.202,1	11,3
Comarques Centrals	Osona	44	138.384	179	0,03%	3.145,1	4,1
Comarques Centrals	Solsonès	6	10.004	118	0,06%	1.667,3	19,7
metropolità de Barcelo	..Baix Llobregat	27	771.516	85	0%	28.574,7	3,1
metropolità de Barcelo	..Barcelonès	9	2.212.658	505	0%	245.850,9	56,1
metropolità de Barcelo	..Maresme	27	413.594	218	0,01%	15.318,3	8,1
metropolità de Barcelo	..Vallès Occidental	21	815.060	39	0%	40.241,0	1,9
metropolità de Barcelo	..Vallès Oriental	40	373.089	102	0,01%	9.327,2	2,6
Penedès	Alt Penedès	27	94.885	72	0,03%	3.514,3	2,7
Penedès	Anoia	23	105.591	69	0,02%	4.590,9	3,0
Penedès	Baix Penedès	14	89.804	58	0,02%	6.414,6	4,1
Penedès	Garraf	6	136.328	76	0%	22.721,3	12,7
<b>Total Barcelona (Administració ...</b>		<b>293</b>	<b>5.397.801</b>	<b>1.805</b>	<b>0,01%</b>	<b>18.422,5</b>	<b>6,2</b>
provincia: Girona							
Comarques gironines	Alt Empordà	62	108.025	367	0,06%	1.742,3	5,9
Comarques gironines	Baix Empordà	32	119.767	286	0,03%	3.742,7	8,9
Comarques gironines	Garroba	19	48.677	163	0,04%	2.561,9	8,6
Comarques gironines	Gironès	25	163.749	77	0,02%	6.550,0	3,1
Comarques gironines	La Selva	23	155.325	336	0,01%	6.753,3	14,6
Comarques gironines	Pla de l'Estany	8	26.061	98	0,03%	3.257,6	12,3
Comarques gironines	Ripollès	19	22.489	167	0,08%	1.183,6	8,8

Il·lustració 38: Los totales a nivel tipo y año del informe Equipaments Comarca

## 6. Conclusiones.

---

En este documento no solo se ha presentado un análisis y diseño del proyecto solicitado, sino que, también muestra como todo el proceso ha sido realizado bajo el paraguas de la diagramación **UML**, con el beneficio de tener toda la definición y diagramación integrada en un único repositorio.

A nivel personal el seguir este método basado principalmente en [ref-LujanS] (pero también de [ref-prat-akoka]) me ha servido y me sirve de guía para afrontar de forma focal un parte del problema, pero a la vez tenerlo dentro de un visión conjunta. Por otra parte la lectura de [ref-Bouman] me ha expuesto las diferentes consideraciones que hay que tener al desarrollar un almacén de datos.

Respecto a la consulta de datos por OLAP al ser este un proyecto de datos cruzados entre varios tipos de estrellas, obliga a que la conexión fuera una consulta SQL y por definición de tablas, lo que se pierde parte del potencial de los esquemas, por otra parte, los datos cruzados, que no dejan de ser hacer combinaciones entre manzanas y peras, hacen que ciertos valores calculados en el cubo sea ilógicos, debido al no poder indicarle a un valor que operación hay que utilizar cuando se agrega en cada una de las dimensiones que le afecta, el sistema solo permite elegir una que aplicará en todas las dimensiones.

## 7. Lineas de evolución futura.

---

Los informes cumplen los requisitos del proyecto y además ofrecer nuevas consultas OLAP, aunque en el enunciado no se habla de gráficos creo que haría falta diseñar uno que resumiera el valor de la información presentada destacando de una manera visual y reveladora la razón de ser de este proyecto.

Para el proceso de transformación, se han creado múltiples reglas de validación, pero no estoy satisfecho del registro de los cambios, en este caso es un proyecto de una sola carga y por lo tanto el control de modificaciones, no es importante, y no he creído necesario utilizar el sistema para *las tablas que cambian lentamente*, **Slowly Changing Dimensions (SCD)**, aunque si me hubiese gustado presentar una sistema de registro de altas y escrituras tal como describí en el modelo lógico del ETL, pero la falta de tiempo me lo ha impedido.

## Glosario

---

**Cubo OLAP** Es una expresión que surge de ver el OLAP como un cubo de varias dimensiones en donde en cada una de las caras hay una dimensión del hecho que se analiza.

**Desnormalizar**, Una vez aplicados los conceptos relaciones en la tablas descerlos para conseguir una optimización de tiempo en las consultas.

**Dimensión**. Una entidad sobre la que se puede medir.

**Estrella**. Es el diseño de crear las tablas de hechos con sus dimensiones, de tal forma que la tabla de hecho solo tiene valores y referencias a las distintas dimensiones.

**Esquemas**. Véase metadatos.

**ETL**. Proceso de extracción, transformación y carga de los datos desde sus distintos formatos y fuentes hasta el almacén de datos.

**Hecho.** ver valor

**HOLAP** Modelo híbrido de almacenamiento OLAP que utiliza base de datos relacional y multidimensional

**Máquina Virtual.** Software que simula un equipo determinado sobre el cual se le instala un sistema operativo funcional

**Mecanismo ETL,** Una unidad elemental de transformación, se ha identificado los siguientes mecanismos: Agregación, conversión, filtro, incorrecto, unión (join), cargador, registrador, fusionar (merge), sustituto, envoltorio.

**Metadatos.** Datos sobre los datos, describen cual es la estructura de los datos que se almacenan y como se relacionan. Incluye esquemas basados en la estructura y el lenguaje del negocio.

**Modelo Conceptual** Resultado de aplicar el análisis centrado en buscar entender el ámbito o dominio real en el que está insertado el problema.

**Modelo Entidad/Relacion.** Modelo en donde se identifican la entidades o elementos que conforman las estructuras de datos y sus relaciones.

**Modelo Lógico.** Es el modelo conceptual transformado para adaptarlo al sistema gestor de la base de datos y teniendo en cuenta los criterios de desnormalización.

**MOLAP** Es un sistema OLAP que su almacenamiento se realizan en base de datos especificas conocidas como multidimensionales.

**OLAP.** (Procesamiento analítico en línea) es la herramienta y las técnicas que permiten realizar consultas decisorias en línea.

**Paquete UML,** Agrupación de varios elementos del dominio del problema o diseño de la solución, se representa como una carpeta.

**ROLAP** Es un sistema OLAP que su almacenamiento se realiza en una base de datos relacional

**Stagging Area.** Espacio de almacenamiento utilizado durante el proceso ETL.

**UML** Lenguaje Modelado Unificado, es una método de desarrollo de proyectos informáticos y un conjunto de técnicas de diagramación

**Valor / Hecho** unidad de información mensurable en varias dimensiones.

## Bibliografía y Enlaces

---

### Enlaces web

- [ref-LujanS] **Sergio Luján-Mora.** *Data Warehouse Design with UML PhD Thesis.* Departament of Software and Computing Systems University of alicante. Junio del 2005. Formato PDF. Enlace web: <http://www.dlsi.ua.es/~slujan/files/thesis.pdf>
- [ref-prat-akoka] **Nicolas Prat\* Jacky Akoka\*\*** *Form UML to Rolap multidimensional databases using a pivot model* Formato PDF. Enlace web: [http://faculty.essec.fr/n.prat/BDA02\\_Pratt\\_Akoka.pdf](http://faculty.essec.fr/n.prat/BDA02_Pratt_Akoka.pdf)
- [ref-argoUML] <http://argouml.tigris.org/> *ArgoUML is the leading open source UML modeling tool and includes support for all standard UML 1.4 diagrams. It runs on any Java platform*

#### Libros

- [ref-Bouman] **Roland Bouman Jos van Dongen** “Pentaho Solutions: Business Intelligence and Data Warehousing with Pentaho and MySQL”. Ed. WILEY.

#### Referencias

- [ref-hüse] **B. Hüsemann, J. Lechtenbörger and G. Vossen**. Conceptual Data Warehouse Modeling. In *Proceedings of the 2<sup>nd</sup> International Workshop on Design and Management of data Warehouses (DMDW'00)*, páginas 6.1-6-11, Stockholm, Sweden, 5 Junio del 2000

## Anexo: Por que **ArgoUML**?

---

Un primer motivo era tener un herramienta de análisis que permitiese definir todo el proyecto y así conseguir tener un repositorio único con todo el análisis/diseño del proyecto,

Un segundo motivo era que quería evitar utilizar varios tipos de herramientas diferentes, con los previsibles problemas de sincronización y manteamiento.

Otros motivos:

- Vengo utilizando **ArgoUML** para mis proyectos profesionales y por la tanto me evita incluir en el **TFC** la planificación de tiempo que requeriría para superar la curva de aprendizaje
- Es una herramienta código abierto, y por lo tanto, colaborativo.

## Anexo Modelo Conceptual basado en UML

---

### **NOTAS SOBRE LA REPRESENTACIÓN DEL MODELO CONCEPTUAL**

Para el modelo conceptual no utilizamos el modelo *Entidad/Relación ER*, si no, aprovechamos de que de hecho el **UML** es un superconjunto de las notaciones **ER** y utilizamos los diagramas de clases con notaciones extendidas, principalmente estereotipos.

En esta ocasión aprovechamos la versatilidad de crear un perfil **UML** específico para los almacenes de datos siguiendo la propuesta de nomenclatura de [ref-LujanS]

Para la representación del modelo conceptual hemos utilizado la propuesta de [ref-LujanS] de dividirla en niveles utilizando la agrupación por paquetes que ofrece el UML, que tiene la ventaja adicional de permitir definir diferentes dominios conceptuales. En la propuesta de [ref-LujanS] indica la utilización de 3 niveles, con la idea de dar cabida a poder diagramar estructuras complejas, y a su vez que éstas se representen con el mínimo número de niveles.

En nuestro proyecto siguiendo la filosofía de mostrar la complejidad de forma entendible y con el menor número de niveles posibles y viendo que la magnitud de nuestro proyecto nos permite quedamos en dos niveles:

- Nivel 1: Definición del modelo, en donde el artefacto **UML paquete** representa el esquema estrella del modelo conceptual. Identificamos con el estereotipo `<<star>>` aquellos que contiene en su interior una estructura de estrella con sus valores y dimensiones (como mínimo a de tener una dimensión) y `<<dimension>>` para aquellos paquetes que en su interior solo contienen dimensiones. Dichas dimensiones serán compartidas por los demás paquetes estrella. Las flechas de dependencia indican que el paquete como minino contiene una dimensión que comparte la estrella
- Nivel 2: Definición de la estrella: Se representan como **clases UML** cual es el hecho ,cuales son sus dimensiones y sus correspondientes niveles de jerarquía. Para el nivel 2 utilizamos los siguientes artefactos para los siguientes conceptos.
  - Hecho o valor : clase con estereotipo `<<fact>>`. Puede tener atributos
  - Dimensión: clase con estereotipo `<<dimension>>`. No podrá tener atributos, ver más



adelante la explicación sobre el uso de base y dimensión.

- Niveles de jerarquías de agregación: clase con estereotipo `<<base>>`, puede tener atributos
- Existe otra clase del tipo hecho que tiene un etiquetado especial `<<degenerateFact>>` para indicar que realmente es una asociación Muchos a Muchos entre un hecho y una dimensión
- Clasificación de las jerarquías: utilizamos el dibujo de la agregación con el estereotipo que indica que permite la navegación `<<rolls-upto>>`
- Las relaciones entre clases pueden tener indicadores de cardinalidad.
- Estereotipo para el atributo que hace funciones de clave primaria: `<<dwkey>>`
- Estereotipo para el atributo que se utilizará como identificador para la navegación por el cubo dimensional en herramientas **OLAP** `<<descriptor>>`
- Estereotipo para el atributo que contiene un valor, `<<factAttribute>>` evidentemente se utilizará en una `<<fact>>`
- Estereotipo para indicar una dimensión que no existe como una entidad independiente, si no, que esta dentro de el hecho `<<degenerateDimension>>`.
- En el nivel conceptual para los tipos de los atributos se utilizan los que están en el perfil de predefinido en **ArgoUML**, *Integer* y *String* para diferenciar se trata de un numero o de un texto, sin entrar en mas detalle.

## CRITERIOS PARA EL USO DE `<<DIMENSION>>` Y `<<BASE>>`

Para la representación de la dimensión y los niveles de jerarquía seguimos la recomendación de [ref-LujanS] de utilizar la definición de Hüsemman [ref-huse] en la que una dimensión contiene una primera jerarquía llamada *dimensión de nivel terminal*, (la representamos con el estereotipo `<<dimension>>`) desde la que se asciende “Rolls-upTo” hasta el nivel de “todos”, que aunque esta última no se representa, implícitamente está en toda dimensión.

La *dimensión de nivel terminal* se utiliza como clase enganche entre estrella y los niveles de jerarquías de clasificación que contiene dicha dimensión. Funciona como “un representante “de la dimensión y quienes definen los niveles de jerarquía son las clases `<<base>>` mediante sus relaciones de tipo agregación entre ellas . Con ello conseguimos 2 objetivos acercarnos al modelo de cubo y por otro facilita la representación del modelo. Para ver de una manera gráfica ir al apartado “*Nivel 2: Paquete dimensión comuns*”

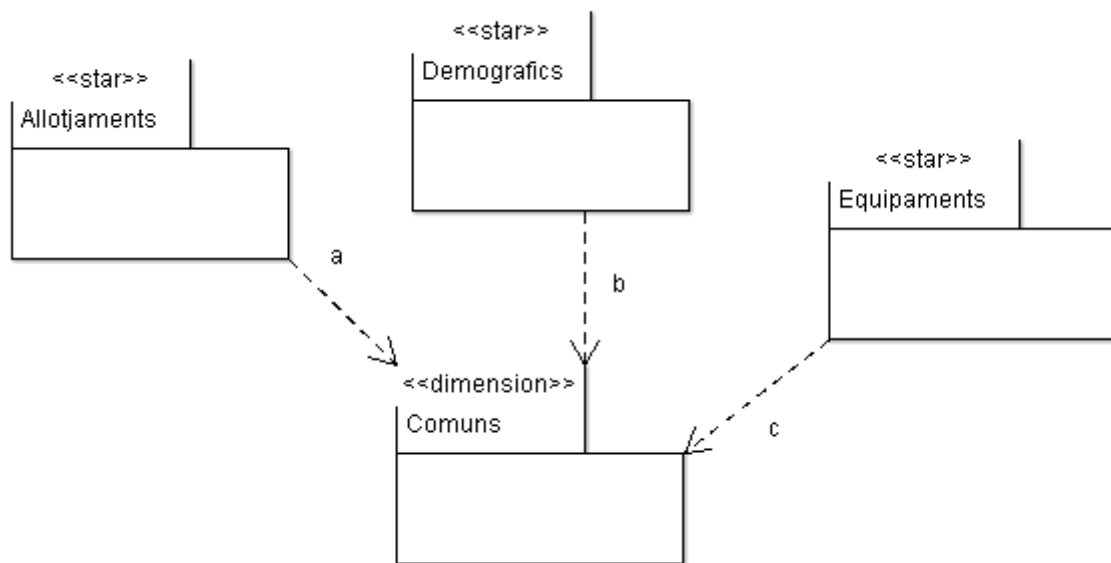
## **NIVEL 1: MODELO CONCEPTUAL**

En el apartado de “Definición del contexto” se han identificado dos grandes grupos de consultas que ha de responder el sistema:

- Comparativas de establecimientos con zona geográficas y su población (habitantes y km<sup>2</sup>), en todos ellos ha de poderse navegar por año y por tipo de establecimiento y agrupación geográfica.
- El sistema a de soportar comparativas y ser capaz de establecer relaciones entre los establecimientos y los equipamientos, Además ha de poderse , navegar por el árbol de tipos de equipamientos y clases de equipamientos.

A mayores hemos propuesto diseñar el sistema para que permita

- Comparativas anuales entre equipamientos.
- Comparativas entre los equipamientos y los datos geográficos y demográficos.



Diagr1: Modelo Conceptual.

A partir de lo anterior expuesto, en nuestra representación tenemos 3 formaciones de estrellas (paquetes <<start>>): *Allotjaments*, *Demografics*, *Equipaments*. Y una dimensión compartida por las tres estrellas llamada *Comuns*. Así tenemos en un primer nivel una visión de los diferentes tipos de valores que conforma el producto.

- El paquete *Allotjaments* contiene la estrella con los valores sobre el alojamiento hotelero y su dimensión de categorías de establecimientos y las dimensiones comunes: *Any* y *Comarca*
- El paquete *Demografics* contiene la estrella con los valores sobre la demografía y extensión geográfica, además de las dimensiones compartidas *Any* y *Municipi*
- El paquete *Equipaments* contiene la estrella con los valores agregados de equipamientos y

su dimensiones de categorías de equipamientos (*Equipament*) e instalaciones (*Instalacions*), además de las dimensiones compartidas *Any* y *Municipi*

## **NIVEL 2: PAQUETE DIMENSIÓN COMUNS**

*Nota: El diagrama se encuentra en la siguiente página y en formato apaisado*

Para la representación de las dimensiones y como ya se indica en el apartado “Criterios para el uso de <<*dimension*>> y <<*base*>>” definimos una que funciona como representante y después están los elementos de la jerarquía de niveles.

Dimensión *Any*:

En este caso tenemos, a la dimensión *Any* que solo tiene el nivel base de *Any*. Este es el -único caso en donde la clave no será asignado automáticamente por la base de datos, sino será el año de cuatro cifras. A su vez será utilizado como descriptor en el sistema OLAP.

Dimensión *Municipi*:

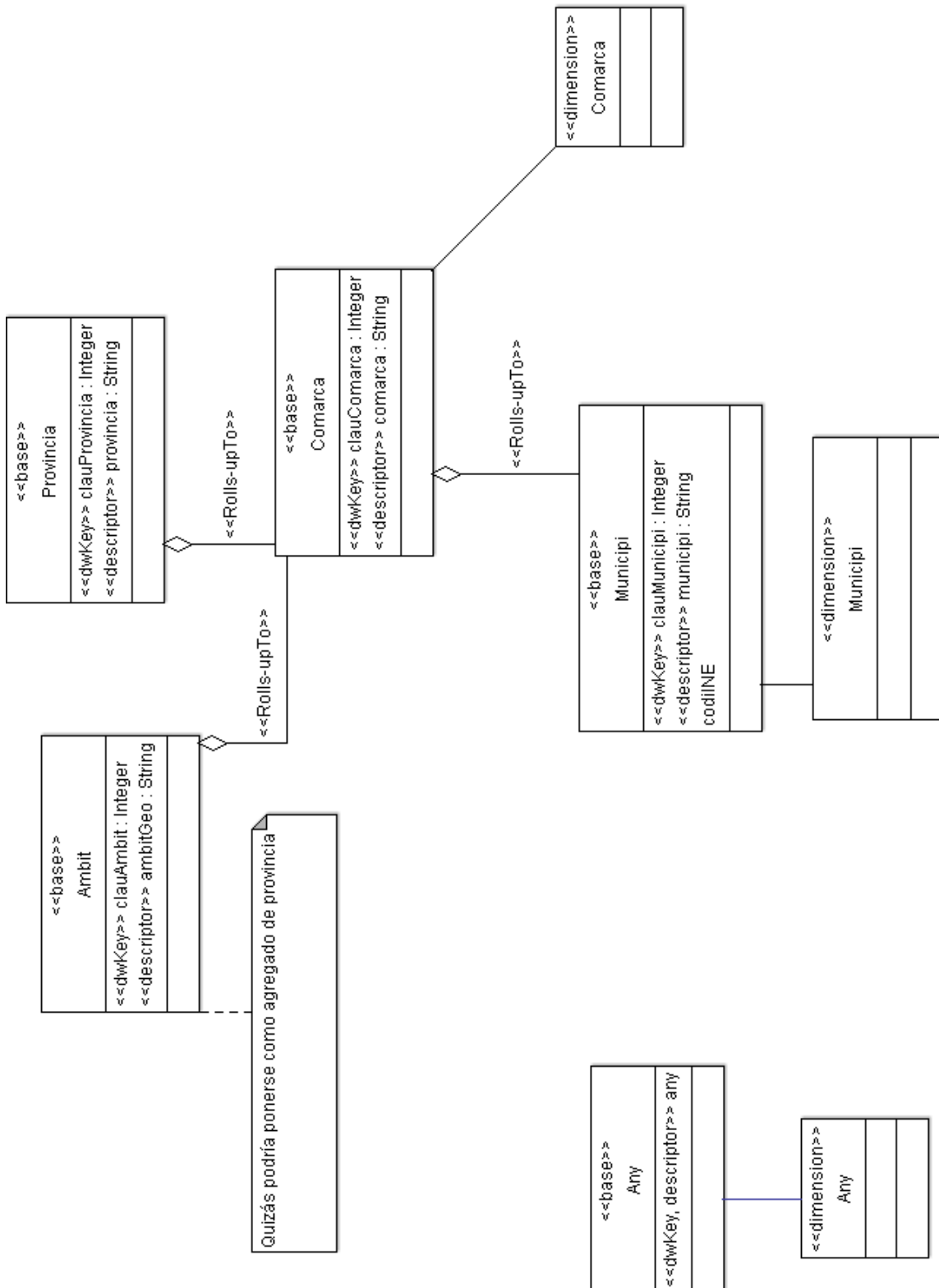
Tiene como base de más bajo nivel al municipio (*Municipi*) que después asciende por *Comarca* que a su vez se puede ascender por *Ambit* y *Provincia*. que como se ve en diagrama comparte con la dimensión *Comarca*

Dimensión *Comarca*

Tiene como base de más bajo nivel a la *Comarca* desde la que se puede ascender por *Ambit* o *Provincia*.

Quizás *Ambit* se podría poner como agregado de provincia, pero como *Ambit* parte de conceptos administrativos diferentes que las provincias, por el momento, se dejan como dos visiones diferentes de agrupar las comarcas.

Comarcas y Municipios: La razón de que existan dos puntos desde donde se “engancha” al árbol de clasificación administrativa, es debido a que existen, *hechos* que tienen su nivel más bajo de representación en el *Municipi* y otros en *Comarca*



Diagr 2: Comuns

## **NIVEL 2: PAQUETE ESTRELLA ALLOTJAMENTS**

*Nota: El diagrama se encuentra en la siguiente página*

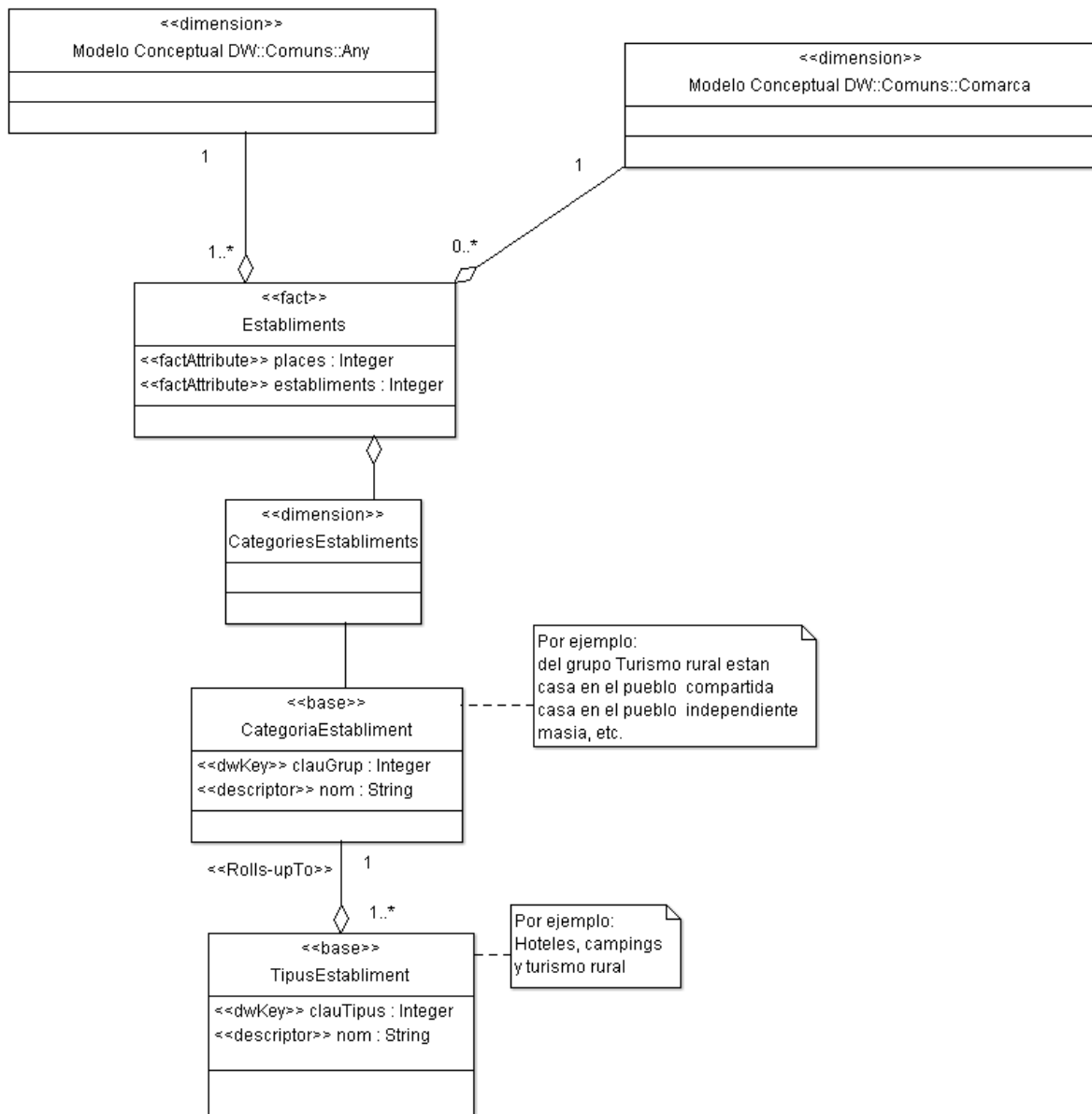
La estrella Allotjaments tiene como el hecho de que exista un establecimiento y su número de plazas. No se puede acumular dichos valores por varios años ya que no tiene sentido.

Las dimensiones tienen como jerarquía de base, el año (*Any*) comarca (*Comarca*) y categoría de establecimiento (*CategoriaEstabliment*) que se puede ascender a los tres tipos básicos (*TipusEstabliment*)

La dimensión *CategoriesEstabliments* tiene la siguiente estructura.

<i>TipusEstabliment</i>	<i>CategoriaEstabliment</i>
<i>Hotels</i>	<i>Hotels (estrellas or)</i>
	<i>Hostals (estrellas argent)</i>
<i>Campings</i>	<i>Luxe</i>
	<i>1</i>
	<i>2</i>
<i>Turisme Rural</i>	<i>3</i>
	<i>Casa de poble compartida.</i>
	<i>Casa de poble independent</i>
	<i>Masia</i>
	<i>Masoveria</i>

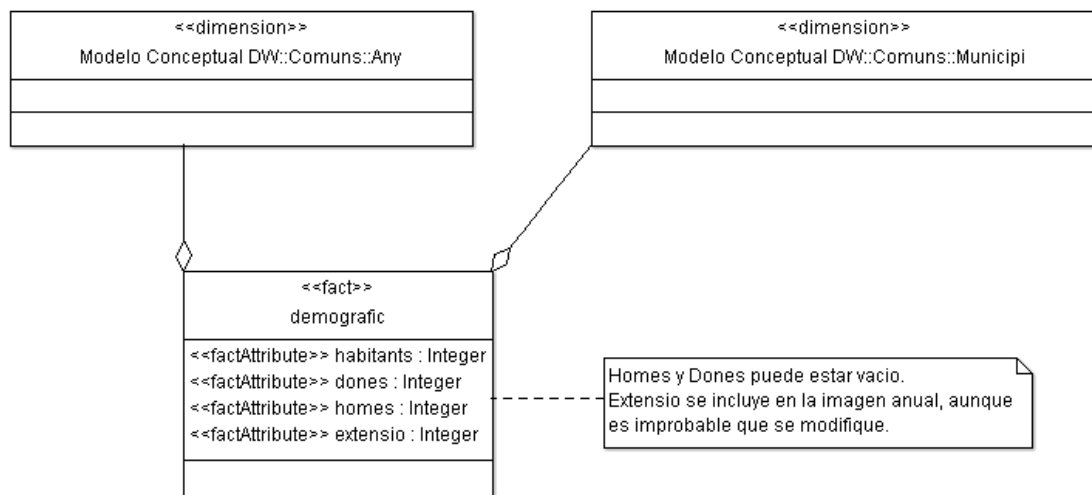
•



Diag3: Paquete estrella *Allotjaments*

## NIVEL 2: PAQUETE ESTRELLA *DEMOGRÁFICS*

La estrella *Demográfics* tiene como el hecho que en una determina zona geográfica viva un individuo, que este sea hombre o mujer y, por otra parte, el dato geográfico de los Km2 cuadrados. En el caso de hombre o mujer las fuentes de datos solo suministran información del año 2012, pero el sistema está preparado para insertarla en otros años. No se puede acumular los valores geográficas por varios años ya que no tiene sentido.



Diagr4 Paquete estrella *Demográfics*

## NIVEL 2: PAQUETE ESTRELLA *EQUIPAMENTS*

*Nota: El diagrama se encuentra en la siguiente página*

La estrella *Equipaments* tiene como hecho la existencia de un determinado tipo de equipamiento y de el tipo de instalaciones que posee. No se puede acumular los valores los equipamientos por varios años ya que no tiene sentido.

En las fuentes de datos la información esta reducida a un solo año, pero el sistema se a diseñado para que permita la entrada de otros años.

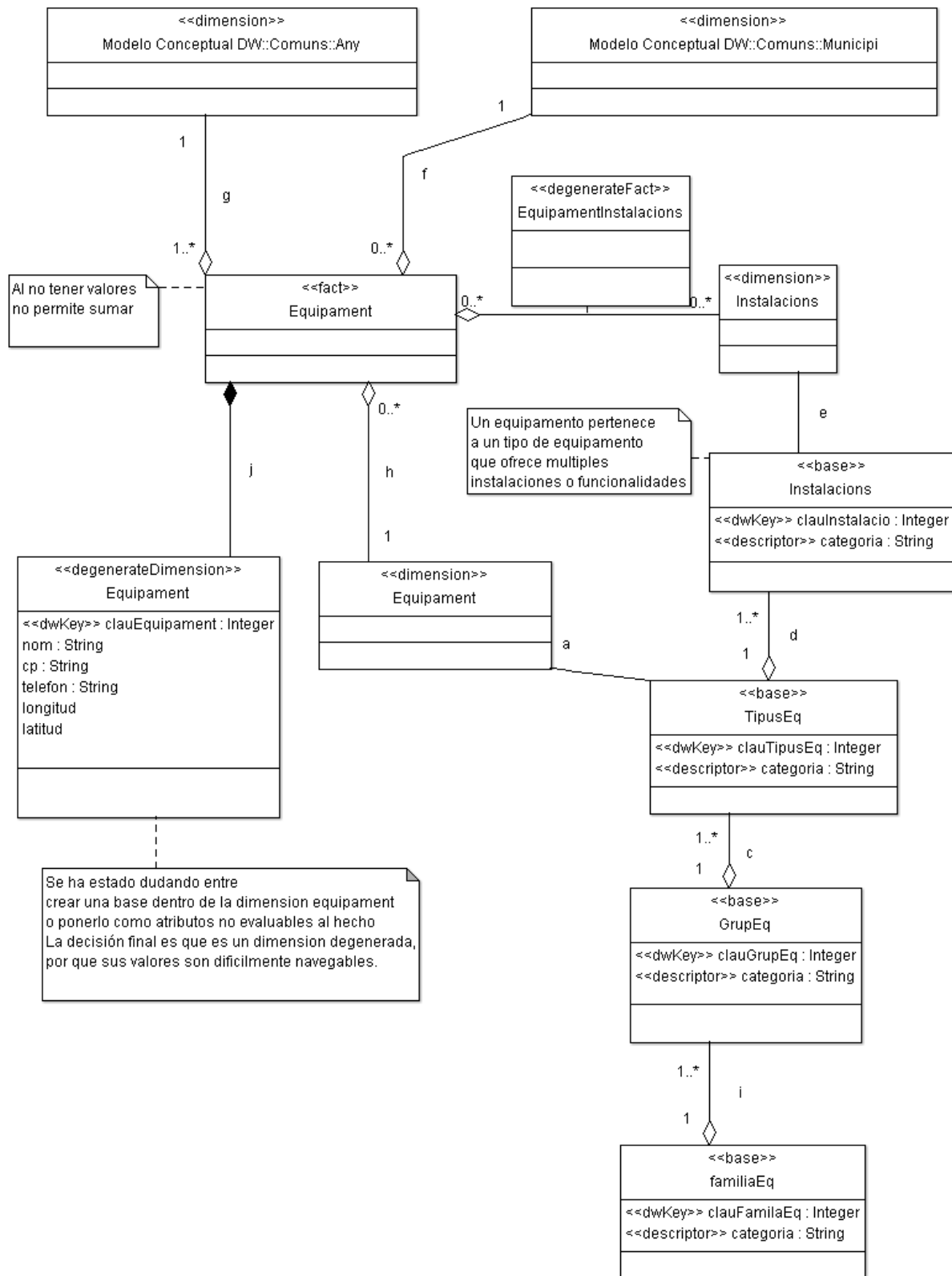
El análisis de las fuentes de datos nos da que la estructura de clasificación de los equipamientos esta formada por 3 niveles jerárquicos que los hemos identificado como *FamiliaEq*, *GrupEq* y *TipusEq*.

A mayores tras estos 3 niveles y en función de *TipusEq* puede el equipamiento contener unas instalaciones para realizar una serie de actividades. El concepto de instalaciones puede hacer referencia tanto a los espacios como a los servicios. Un ejemplo del segundo caso son los centros educativos que pueden dar unos o otros niveles de enseñanza.

En este diagrama existen varias decisiones sobre la definición del dominio a destacar:

- Como se analiza el hecho de que exista un equipamiento de un tipo y con unas instalaciones sea un valor, la información que acompaña al equipamiento se puede decir que es adjunta al hecho, y por lo tanto, se entiende que no se utilizará para navegar por el cubo multidimensional así que se ha agrupado en una dimensión del tipo degenerada
- Como la lista de instalaciones (*Instalacions*) que tiene el hecho equipamiento (*Equipament*) es múltiple y un instalacion puede estar en múltiples equipamientos existe el patrón de que un hecho *Equipament* con respecto a una dimensión *Instalacions* tiene una relación N a M y por lo tanto hay que indicar que es una hecho degenerado.





Diagr5: Paquete estrella *Equipaments*

## Anexo Modelo Lógico basado en UML

---

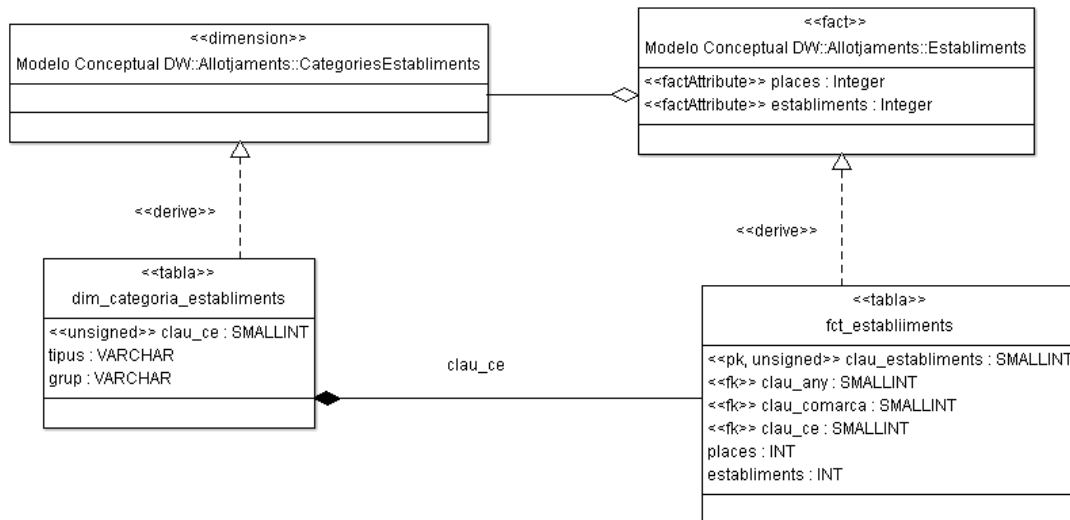
### **NOTAS SOBRE LA REPRESENTACIÓN DEL MODELO LÓGICO**

Para el modelo lógico también se ha creado un perfil específico para poder representar el modelo de base de datos. Así tenemos:

- `<<derive>>` Se utiliza conjuntamente con la flecha de tipo realización (en **OO UML** es utilizada para marcar la relación entre una interfaz y su implementación) en este caso es para indicar la relación entre el elemento del modelo conceptual y el modelo lógico. Apunta a que dimensión o hecho se deriva la tabla,
- `<<tabla>>` a nivel de clase para indicar que es una tabla.
- `<<fk>>` a nivel atributo indica que es una clave foránea.
- `<<pk>>` a nivel atributo indica que es una clave primaria
- `<<unsigned>>` hace referencia a que es un tipo numérico sin signo.
- Están definidos los tipos de atributos más comunes en base de datos y utilizados por **MySQL**: BIGINT, BINARY, BLOB, BOOL, CHAR, DATE, DATETIME, DECIMAL, DOUBLE, ENUM, FLOAT, INT, MEDIUMINT, SMALLINT, TEXT, TIME, TIMESTAMP, VARBINARY, VARCHAR, YEAR.

Expongo en la página siguiente una parte del diagrama del modelo lógico para *Allotjaments* para destacar dos tipos de conexiones que pueden parecer contradictorias entre el modelo conceptual y el modelo lógico.

- En el modelo conceptual las dimensiones son agregados de los hechos.
- En el modelo lógico las tablas de hecho son composición de las dimensiones.



Diagr6 Contradicción en las relaciones.

Esto es resultado de la diferente visión que se aplica en ambos casos, en el modelo conceptual se muestra la centralidad del diseño estrella y por lo tanto, las dimensiones son componentes del hecho o valor, pero en el diseño lógico, se muestra que no puede existir un valor si no existen las correspondientes tablas de dimensiones. Así lo marcan las claves foráneas y queda reforzado por la relación de tipo composición

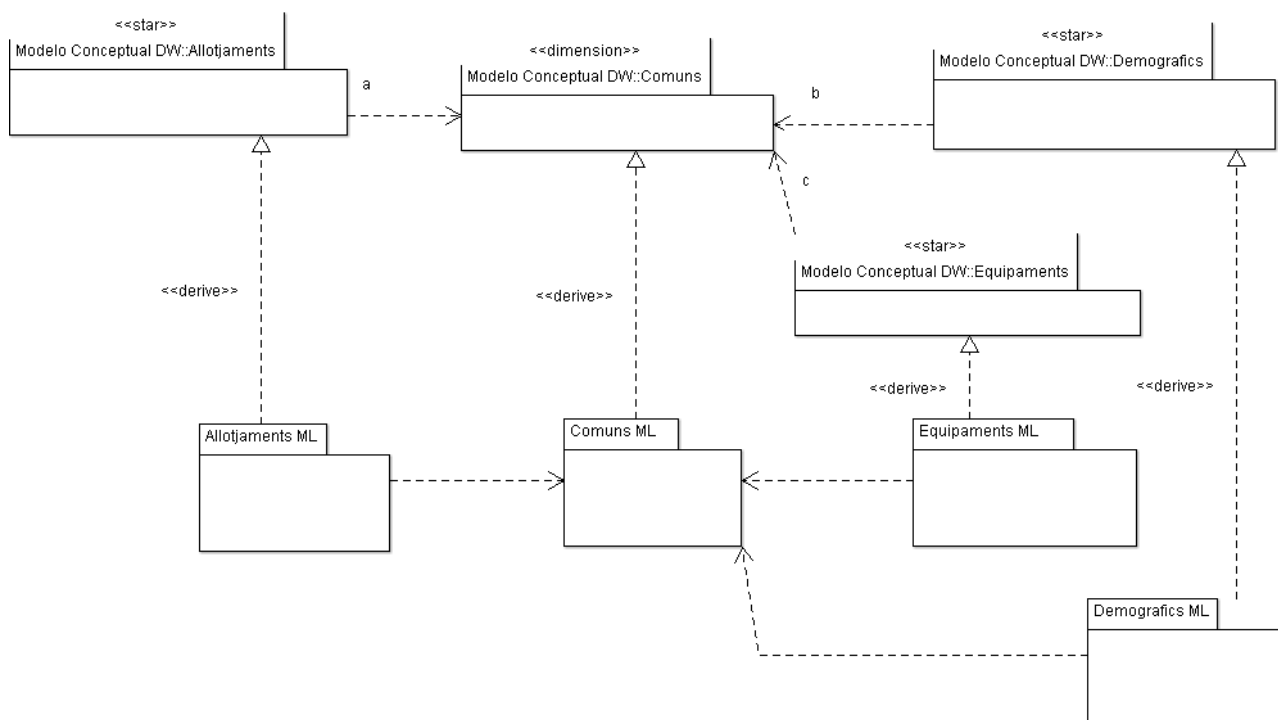
## NOMENCLATURA DE LOS ELEMENTOS

Definimos en el modelo lógico los nombres finales que poseerán los diversas tablas y campos, por lo que se han establecido una serie de convenciones, siguiendo un versión libre del modelo propuesto en [ref-Bouman] .

- prefijo **dim** indica que la tabla es una dimensión.
- prefijo **fact** indica que la tabla es un hecho o valor
- prefijo **agg** indica que es tabla con valores agregados de otras tablas.
- prefijo **deg** indica que es un tabla resultado de una “*dimensión degenerada*”
- En atributos:
- prefijo **clau** indica que es una columna que forma parte de la clave.

## MODELO LÓGICO NIVEL 1: VISIÓN GENERAL.

En el diagrama de la visión general del modelo lógico del almacén de datos, nos permite ver la derivación (relación) que existe con el modelo conceptual. Para buscar ser homogéneos se han aplicado los mismos criterios en la división en paquetes. Con ello, además, facilitamos el mostrar la relación entre los modelos.



diagr7: Modelo Lógico Multidimensional

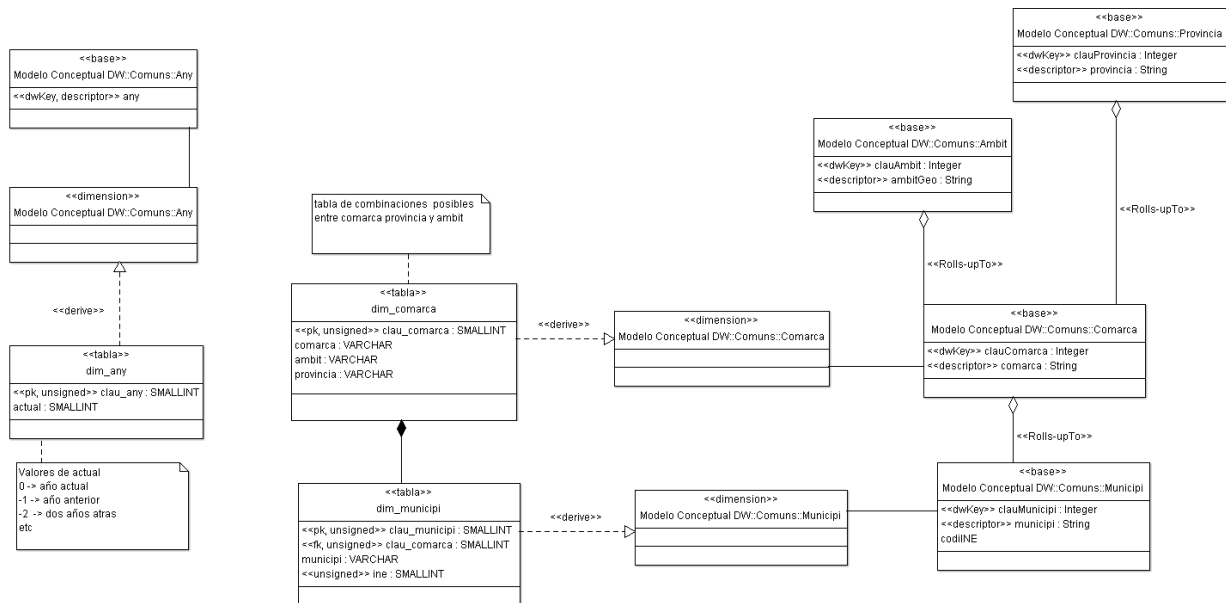
## NIVEL 2: PAQUETE LÓGICO COMÚNS.

Nota: El diagrama se encuentra en la siguiente página y en apaisado.

En la diseño de la tabla del año (*dim\_any*) se ha incluido el nuevo atributo *actual* que sirve para indicar cual es el año actual y cuales son los años anteriores, esto, por ejemplo, aporta rapidez en las selecciones del tipo comparativa entre este año y los 3 anteriores.

En el diseño de la tabla de comarca (*dim\_comarca*) se ha decido desnormalizar y agrupar todos los elementos en un sola tabla y que en ella estén todas las combinaciones posibles. En esté caso al ser jerárquico solo hay una combinación por *comarca*, *ambit* y *provincia*.

Las claves de *dim\_comaca* y *dim\_muicipi* serán autogeneradas por el sistema gestor de bases de datos.

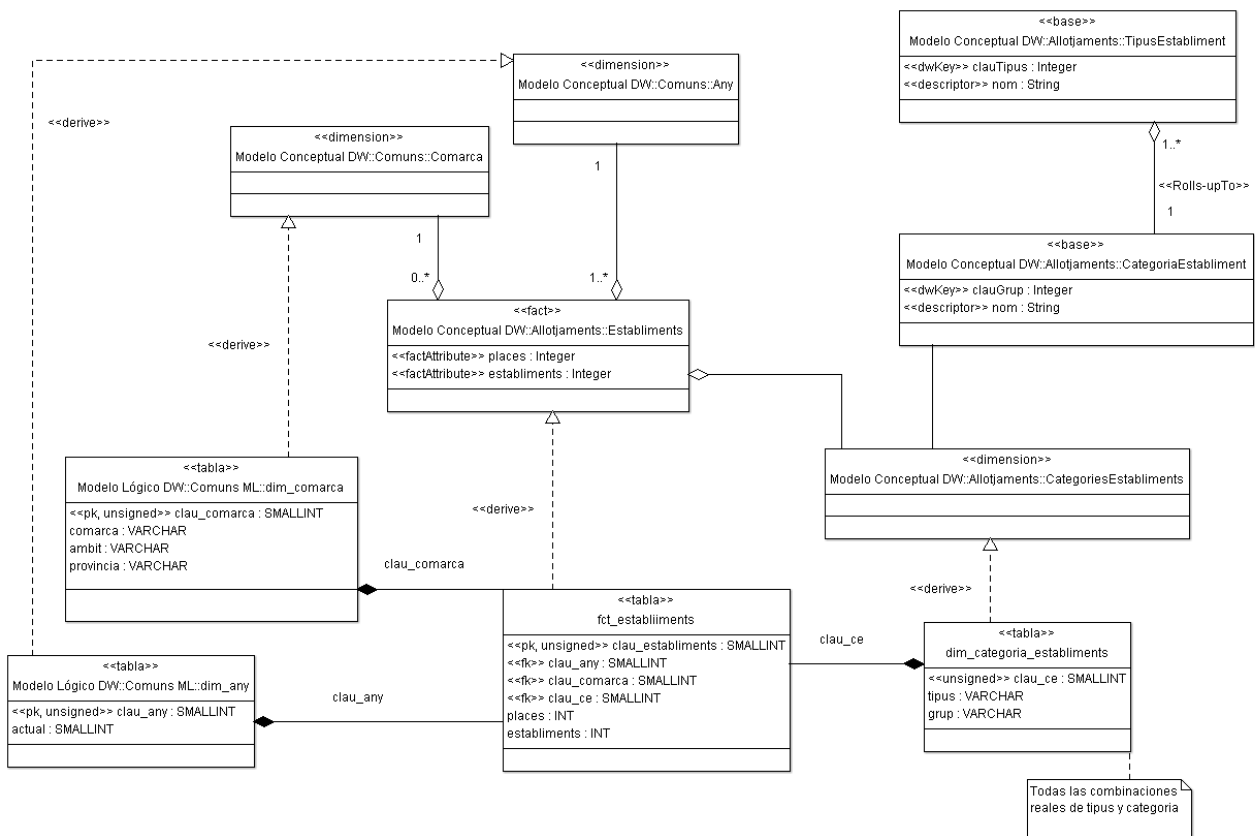


diagr8: Paquete Lógico común

## NIVEL 2: PAQUETE LÓGICO ALLOTJAMNETS.

En el diseño de la tabla de categorías (*dim\_categoria\_establiments*) se ha decido desnormalizar y agrupar todos los elementos en un sola tabla y que en ella estén todas las combinaciones posibles. En esté caso al ser jerárquico solo hay una combinación por *tipus* y *grup*.

Todas las claves primarias de las tablas sobre alojamientos serán generadas por el sistema gestor de base de datos.

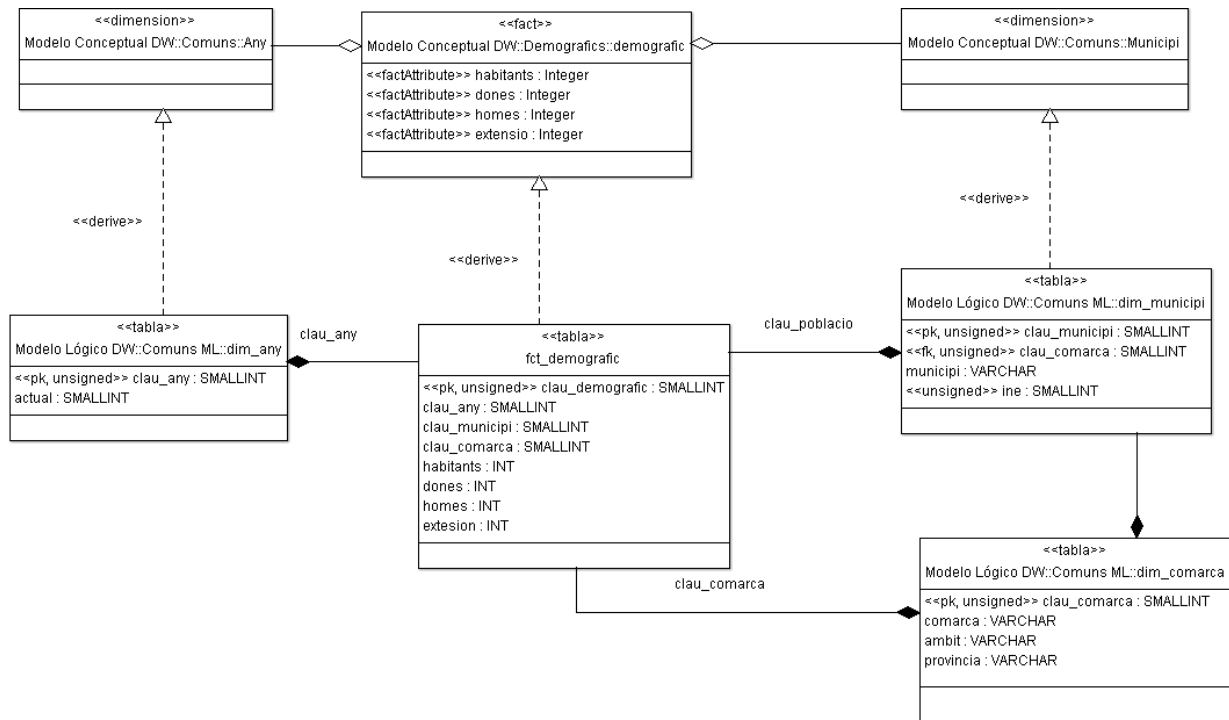


Diagr9

## NIVEL 2: PAQUETE LÓGICO DEMOGRAFICS.

La clave primarias *fct\_demografic* será generada por el sistema gestor de base de datos.

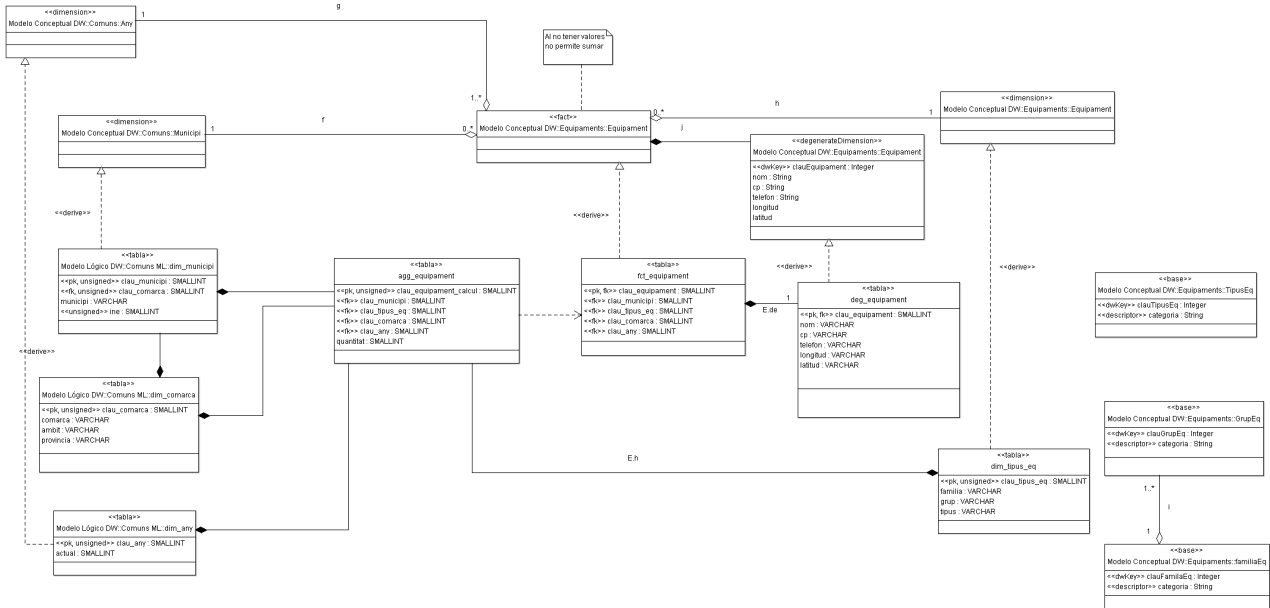
Los atributos Dones y homes pueden tener valores nulos. Es importante entonces que los informes sepan si se va consultar ese dato en un año que no tenga dichos valores.



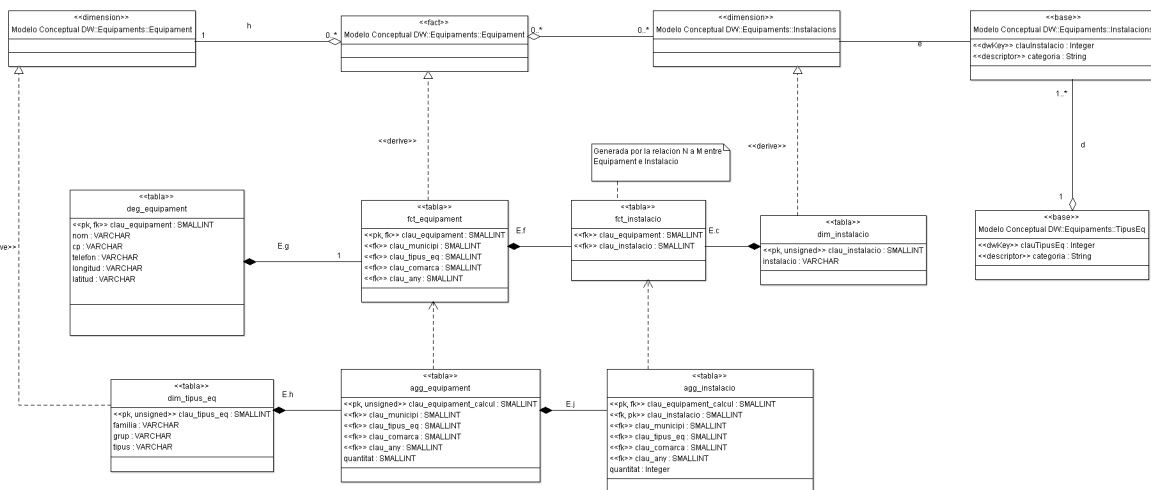
Diagr10

## NIVEL 2: PAQUETE LÓGICO EQUIPAMENTS.

Como lo que es de información relevante es el número de establecimientos de un tipo y sus instalaciones, se han definido tablas agregadas que contienen dicho valor ya totalizado al nivel más bajo de la jerarquía.



Diagr11



Diagr12



## Anexo Modelo Lógico ETL basado en UML

---

Para definir el proceso de extracción, transformación y carga de las fuentes al almacén de datos se ha creado un modelo lógico, en el que mediante una representación por niveles se han detallado los diferentes elementos que intervienen en dicho proceso. Dicho modelo se ha basado en las propuestas de [ref-LujanS] más las recomendación para un buen proceso de **ETL** de [ref-Bouman]

### **MODELO LÓGICO DEL PROCESO ETL**

Se ha decidido crear la definición del proceso **ETL** utilizando el modelado **UML**, con ello se consigue que en un único repositorio de información estén tanto la información sobre el almacén de datos como los procesos de carga y lo que es más importante como se relacionan entre ellos.

Otro aspecto no menos importante es que se define y se tipifican todos los elementos implicados en un **ETL**, con lo cual, se consigue por un lado normalizar la nomenclatura utilizada y por otra facilitar el análisis y diseño ya que esté se centra en ver como se relacionan y se utilizan dichos componentes para que una fuente de datos concreta se traslade al almacén de datos definido.

El paradigma que persigue el modelo se puede resumir en los siguientes puntos (basados en [ref-LujanS] y [ref-Bouman])

- El proceso será descompuesto en unidades elementales llamados **mecanismos ETL**
- Los datos de cualquier fuente que no estén en tablas será previamente traspasado a una estructura de base de datos relacional.
- Existirá una zona en la que se almacenarán todas las tablas necesarias en el proceso de transformación,
- El proceso de transformación, será auditado y en el caso de errores en los datos o en su formato el flujo se detiene en el estado en que se ha producido el error.
- Para los fuentes de orígenes no relacionales se realizara un primer paso que será tener dichas fuente en un formato relacional y sin ningún tipo de manipulación. Es importante que los datos y su estructura no sean manipulados, eso se deja a los **mecanismos ETL** de transformación.
- Las tablas intermedias generadas durante el proceso de transformación se inicializarán en cada proceso de importación.

A continuación desglosamos las unidades elementales llamadas **mecanismos ETL** definidos por [ref-LujanS],

- **Agregación (*Aggregation*)**, como su nombre indica se refiere cuando se totalizan o calculan datos basados en algún criterio
- **Conversión (*Conversion*)**: Cambia los tipos de datos y los formatea o genera nuevos datos a partir de datos existentes.
- **Filtro (*Filter*)**: filtra y verifica los datos.
- **Incorrecto (*Incorrect*)**: Redirige (almacena) datos incorrectos.
- **Union (*Join*)** Se unen dos o mas fuentes en una nueva que sea la combinación de algunos de sus atributos.

- **Cargador (*loader*)** Realiza una carga de datos al almacén de datos.
- **Registrador (*log*)** Registra la actividad de los demás mecanismos.
- **Fusionar (*Merge*)** Integra dos o mas fuentes de datos con atributos compatibles
- **Sustituto (*Surrogate*)** Genera una clave primaria sustituta a la que sería la clave natural que define a la entidad. Por ejemplo agrupar los atributos en una clave natural en un única clave en formato numérico.
- **Envoltorio (*wrapper*)** Extrae los datos una fuente de datos externa para insertarlo en los registro de la base de datos del proceso de **ETL**

Puede pensarse que los mecanismos **ETL** son procesos y no estructuras de datos, está sería un visión limitada, realmente hay que verlo como estructuras de información que han sufrido una manipulación o selección, es decir, reflejan un estado en el proceso de transformación. Por lo tanto, con toda probabilidad, en la implementación tendrá una tabla que contendrá los datos en dicho estado, es decir, tras aplicar las acciones del mecanismo **ETL**. El único que escapa a dicha visión podría ser el mecanismo cargador.

## **NOTAS SOBRE LA REPRESENTACIÓN DEL MODELO ETL**

Para la representación se utiliza el diagrama de clases, en donde los mecanismos ETL son clases identificadas con estereotipos. Las operaciones que se realizan sobre los mecanismos ETL estarán comentadas como notas a la entidad. Se utilizará las conexiones del modelo de clases para representar las relaciones entre los diferentes mecanismos.

Para el modelo **ETL** también se ha creado un perfil específico con todos los estereotipos que intervienen en el proceso **ETL**: <<*etl.agregation*>>, <<*etl.convesion*>>, <<*etl.filter*>>, <<*etl.incorrect*>>, <<*etl.join*>>, <<*etl.loader*>>, <<*etl.log*>>, <<*etl.merge*>>, <<*etl.surrogate*>>, <<*etl.wrapper*>>

Se ha añadido: otros indicativos para diferenciar fuentes de datos: <<*etl.csv*>>, <<*etl.txt*>> y otro para indicar si es importante remarcar que es un proceso manual <<*etl.manual*>>

Se utilizaran las flechas de dependencia para indicar de que otros mecanismos y tablas depende un mecanismo o a que tabla del modelo multidimensional van a actualizar.

## **CRITERIOS GENERALES PARA LOS PROCESOS ETL**

Hay que ir con cuidado con la importación ya que las fuentes proveen valores durante de varios años, y en algunas de estos años los nombres identificativos varían, así como las categorías que incluyen.

Como criterio de diseño se ha establecido que siempre que sea posible utilizar la definición de nombres y las clasificaciones que viene con las fuentes del año 2012.

Todo el proceso de ETL se desarrollará en una zona diferente al almacén de datos.

Los fuentes que vengan exportados de otros sistemas, en la primera fase será importarlos sin cambios en una tabla mimética a la fuente.. Las transformaciones posteriores se realizan en nuevas tablas, dejando la original (la tabla mimética) sin cambios.

Para llevar una auditoria de cuando se produce una ETL en una tabla se ha centralizado en un

tabla *aud\_tables* en donde se indicara la fecha, el número de versión y si es una inserción, una modificación o un borrado (lógico ya que no se producen borrados de registros)

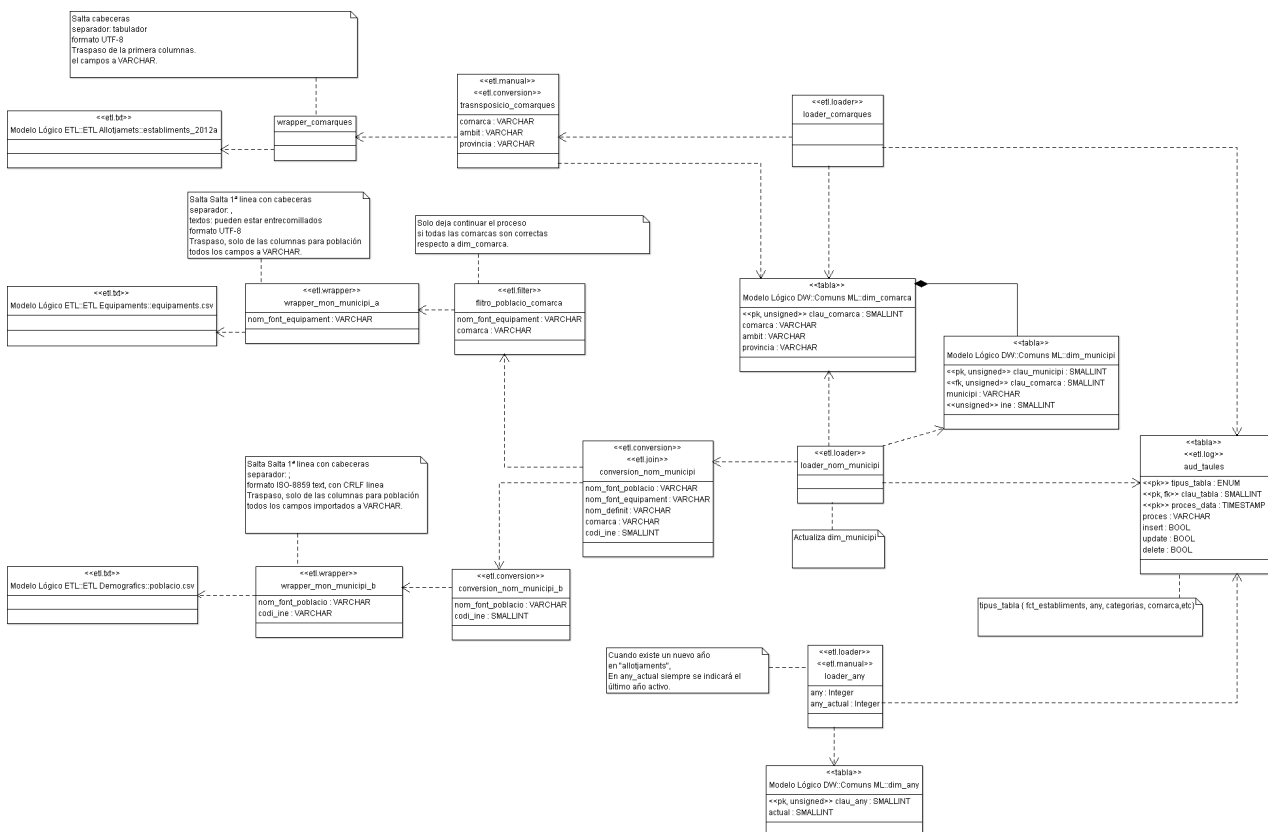
## NIVEL 2: PAQUETE ETL COMUNS.

La carga de la tabla *dim\_Any* se produce cuando se decide importar un nuevo fichero csv de un año en concreto. Se indica que es un proceso manual, ya que se entraran los datos del *loader\_any* manualmente para después realizar el proceso de carga, que puede implicar el cambio de los valores del año actual.

Para determinar la lista con los nombres normalizados del municipio se basará en la información que suministran las fuentes de *poblacio.csv* y *equipaments.csv*.

Para determinar la lista de comarcas, ámbitos y provincias se utilizará la información que suministra la fuente de datos *establiments\_2012* y con ella se normalizarán los nombres, para futuras importaciones.

La tabla del *conversion\_nom\_municipi* como es una tabla que contiene información de la diferentes formas en que se nombra a un municipio quizás sea interesante que se mantenga para futuras importaciones y no se inicie al final o al principio de un proceso.



Diagr4

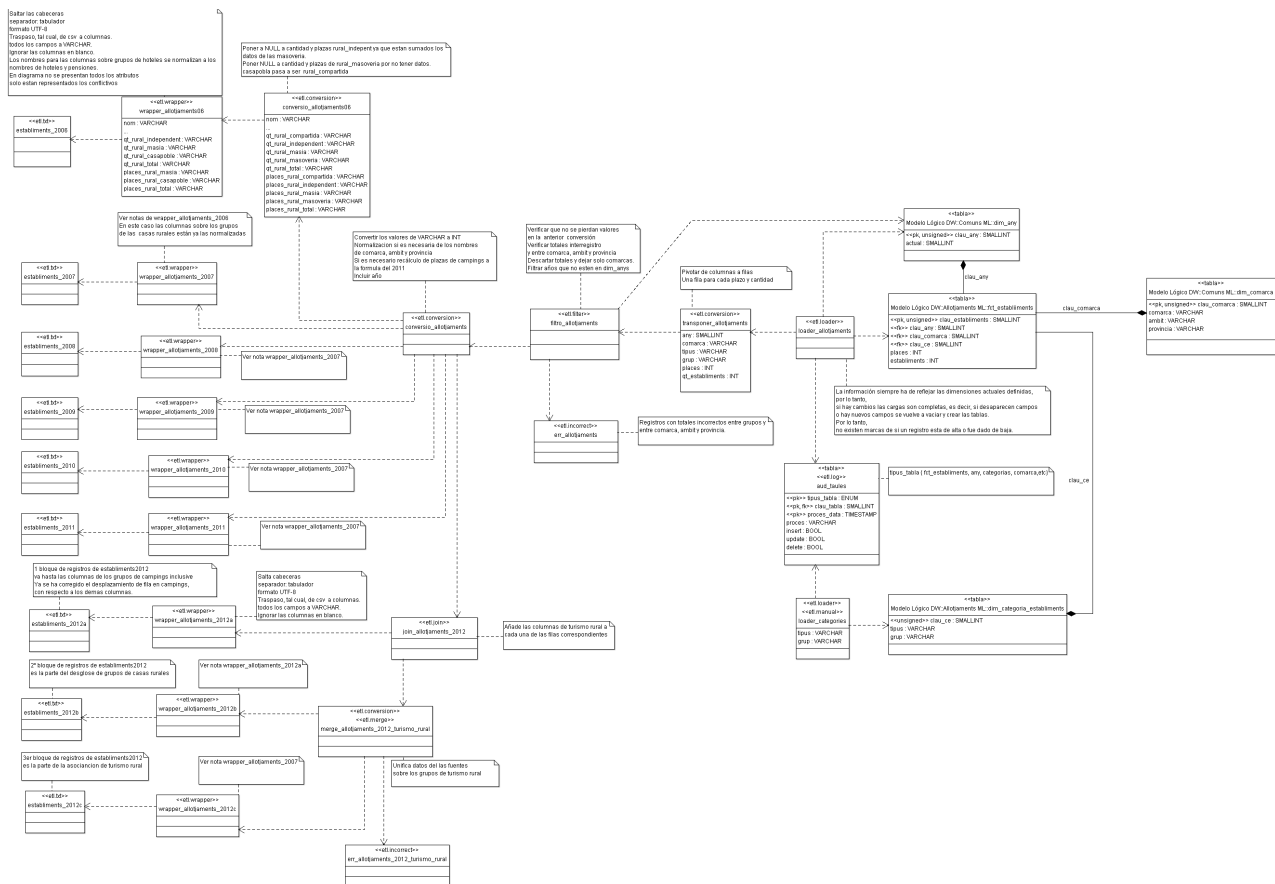
## NIVEL 2: PAQUETE ETL ALLOTJAMENTS.

Este es un diagrama más complejo ya que se analiza cada uno de los ficheros planos con los datos de cada año, haciendo indicaciones precisas de que ajustes hay que realizar en cada uno de ellos.

En el caso del 2012 se ha dividido el fichero original en otros 3 ficheros en que se guardan cada uno de los tres grandes bloques que contiene.

Los datos para el mecanismo *loader\_categories* son entrados manualmente y surgen a partir de las cabeceras del año 2012L.

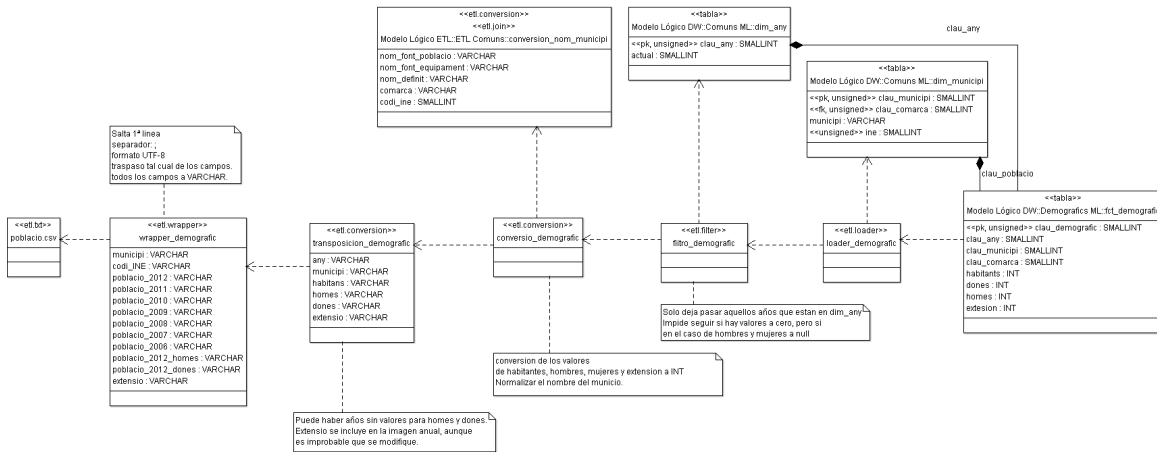
Para importación de los próximos años se aconseja conseguir una versión más próxima a un fichero csv, con lo que se facilitará mucho el trabajo.



Diagr15

## NIVEL 2: PAQUETE *ETL DEMOGRAFICS*.

No hay información adicional que exponer sobre el diagrama presentado. Solo indicar que cuando venga futuros ficheros habrá que revisar si vuelven a repetir los últimos 6 años o viene solo el nuevo año.



Diagr16

