



binary team

عدد الصفحات : 6

م.عبد البديع مراد

المحاضرة : 5

قواعد المعطيات المتقدمة

تعرفنا سابقاً على Model تمثل بـ Relational Database وهي ERD وتعتمد العلاقات One-To-Many , Many-To-Many سنتعرف في هذا المحاضرة على الـ Cube وأبعاده فهو عبارة عن Multi Dimensions Model وله نوعين :

1. Star Schema : لا تأخذ بعين الاعتبار الـ Normalization وهو الأكثر استخداماً .
2. Snowflake Schema : تأخذ الـ Normalization بعين الاعتبار .

تذكرة : تعني الـ Normalization تنظيم الأعمدة والجداول وتبسيط التصميم في قاعدة المعطيات .

كلا النوعين السابقين يعتمد بشكل أساسي على الـ Dimension Table والـ Fact Table :

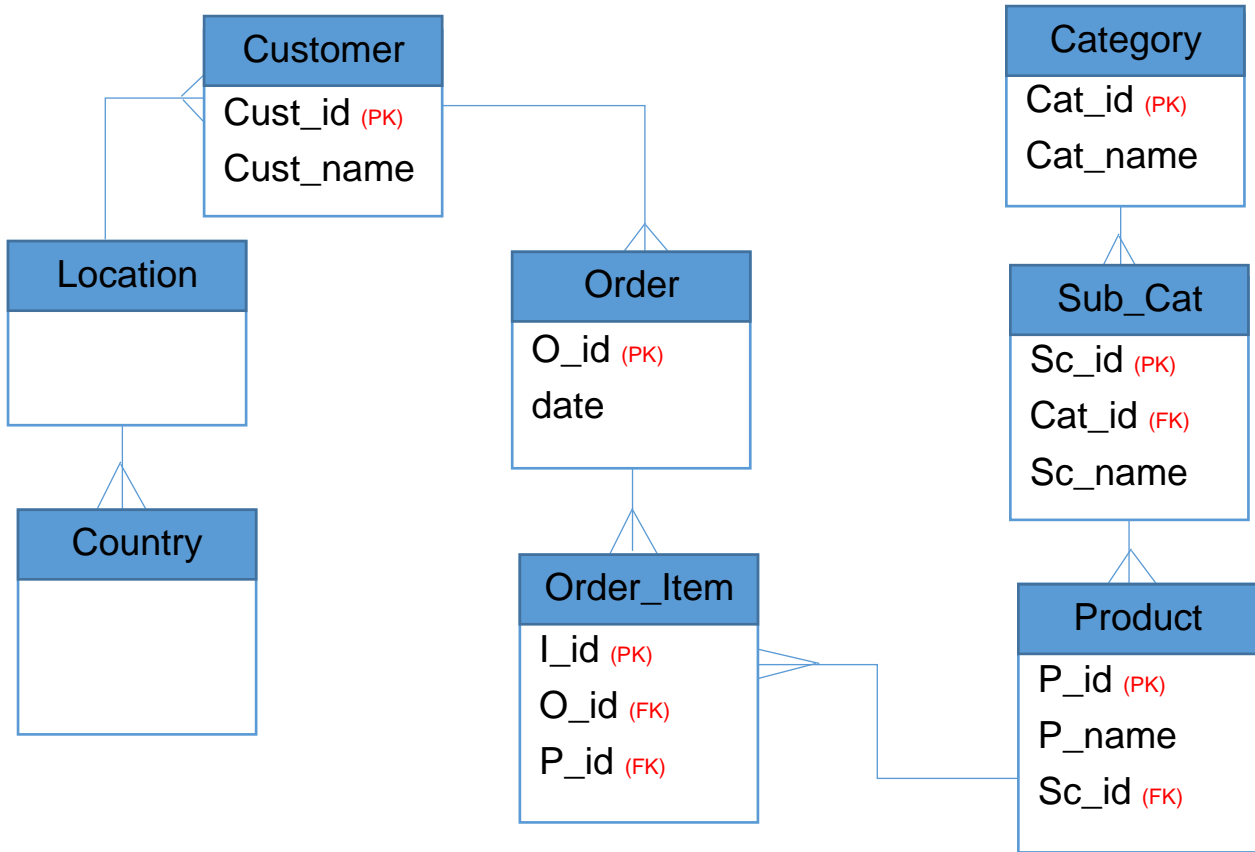
- Dimension Table : يعين الأبعاد ويهم فيه الـ Primary key والـ Hierarchy التي تسمح بعملية Rollup و Break down .
- Fact Table : يهمني الـ foreign key والـ measure .

الـ Dimension Model هو Entity Relationship Model ولكن بأداء أفضل عن طريق وجود بعض التغييرات البسيطة في المفهوم .

الانتقال من ERD Model إلى Dimensions:

- تقريباً 70% مما سنقوم بفعله يحدده الزبون , أي المعلومات الإحصائية التي يريد أخذها من الـ DB بالتالي كي نقوم بعملية الانتقال يجب مراعاة متطلباته .
- بعد الانتقال نقوم بعملية ترحيل المعلومات كل ستة أشهر أو سنة .
- سنأخذ مثال لفهم الانتقال .

مثال : لنحول المخطط التالي إلى Dimension Model



نعمند في عملية التحويل من Relational إلى Dimension بنسبة 30% إلى 40% على أفكار Technical ومن 60% إلى 70% على متطلبات الزبون كتنقيس وأبعاد .

مثلاً الزبون يهيمه أن يعرف كمية مبيعاته ضمن المنتج ومن هو الـ Customer وإلى أي دولة وما هي الـ Category بالتالي قام بتحديد ما يريده فنقوم بعملية التحويل بمراعات متطلباته.

كي نحقق التحويل يجب أن نفكر ما هي الـ Fact Table مع الـ measures وما هي الأبعاد الأساسية Dimensions فنستطيع حينها تحديد الجداول ومحتوى كل منها .

- الـ Fact Table : (هو الـ Root والقلب الأساسي) ينشأ من الجداول التي يجري عليها أكثر Transaction في الـ DB وفي مثالنا هم Order والـ Order_Item ويمكن دمجهما في Fact Table واحد أو اثنين مختلفين ويعود هذا الشيء حسب ما يريده الزبون من متطلبات ودقة بالمعلومة.

السطر في الـ Multi Dimension تقابله n Row في الـ Relational (يقابله سطر أو أكثر من سطر) ومن المستحيل أن يوجد سطر وفي الـ MD ليس له مقابل في الـ Relational .

نستخدم في الـ Fact Table الـ measures وهي aggregation (Count , Avg , Sum) أي قيم تجميعية.

- الأبعاد Dimensions: جداول تنتج عن الـ Tables المرتبطة ارتباط وثيق ومباشرة بالـ Transaction tables والتي هي في مثالنا Order و Order_Item فتشكل الجداول المرتبطة بهما الـ Dimensions ، في هذا المثال تكون الأبعاد هي: Customer , Product وهناك بُعد (موجود دائماً بشكل تقريبي) مهم جداً وهو الـ Date .

مثلاً : أريد معرفة كمية المبيعات كعدد و كـ total price (هاتين القيمتين هما measures).

فيكون لدي Row فيه ال Total Price هو 1,000,000 هل هذا المبلغ تحقق من بيعة واحدة أو أكثر من بيعة ؟ هل من Customer واحد أم أكثر ؟؟ هذا ما تحدده الأبعاد .

ماهي ال hierarchy ؟

معلومات تعبر عن هرمية التركيب وتأتي من الجداول المرتبطة بجداول ال Dimensions .
مثلاً أريد معرفة مبيعات هذا المنتج لكن أريد أن اعرف هذا المنتج من أي Subcategory ومن أي Category

Category → SubCategory → Product

جواكيت رجالي /نسائي البسة

تذكرة: الفرق الأساسي بين ال Star Schema وال Snowflake Schema هو أن الأولى لا تطبق مبدأ ال normalization.

: Star Schema

نقوم بتحويل المثال السابق إلى ال Dimension Model من نوع Star Schema .

1- Fact Table : نسمة Fact Order وينشأ من الجداول Order , Order_Item

وال measures هنا ممكن أن يكون quantity وال total price حيث نريد قياس كمية المبيعات والسعر الإجمالي .

2- Dimension Table : Date , Customer , Product وتكون علاقة ال Dimensions مع ال Fact Table علاقة

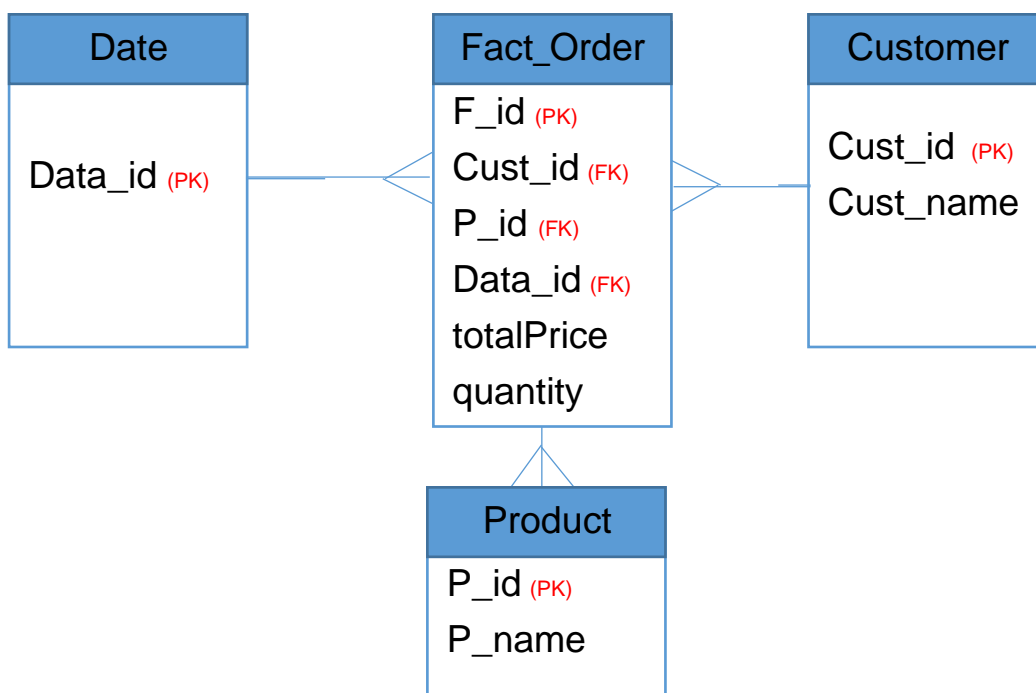
One-To-Many بحيث يوجد في ال Fact Table مجموعة foreign key تعبر عن ال Dimensions .

ال row في ال Customer وفي ال Product

| | |
|---|----|
| 1 | p1 |
| 2 | p2 |

وفي ال date لدي أيام.

| | |
|---|------|
| 1 | أحمد |
| 2 | سمير |



ويكون مثلاً شكل السطر في ال Fact_Order :

| Cust_id | P_id | Date | Total price | Quantity |
|---------|------|----------|-------------|----------|
| 1 | 5 | 2017-1-1 | 100 | 50 |

هذا يعني أن الزبون 1 اشترى المنتج 5 في اليوم 2017-1-1 و كان ال total price 100 والكمية 50 وهذا ال row ممكن أن يكون قد تشكل نتيجة 50 row في ال ERD الأساسي أي هناك 50 Order أو مثلاً نتيجة 3 Order وكل Order فيه quantity مختلفة . أصبح هذا ال row جزء من ال Cube بحيث لدي بُعد لل Product وال Date وال Customer ما يعني أن أي نقطة من هذا الفضاء هي عبارة عن Row ولهذه النقطة عدة أبعاد في مثالنا . من المستحيل أن يكون هناك record في ال Fact_Order ليس لديه أب في ERD لأن عملية ترحيل ال Data من ال ERD إلى ال MD يقوم بها ETL (Extract Transformation Load) وهي أداة جاهزة . لا ننصح بزيادة ال Dimensions لكي لا نصعب الأمور فإذا كانت ال Dimensions كثيرة نبني مكعبين يحوي كل منهما 3 – 4 Dimensions وندمج بينهما بال Tool . عند التحديث عن Cube لا تكون ال query مثلاً select from table وإنما select from cube وهي لغة query أخرى تسمى MDX (Multidimensional Expressions) لن نتعلمها في منهاجنا بل سنتعلم ال structure فقط . نلاحظ أن شكل الجدوال وترباطها يشبه شكل النجمة.

: Star Schema and Hierarchy

نلاحظ أن ال Product له أب وهو Sub category والذي له أب Category .

نفس الأمر بالنسبة لل Customer له أب Location وله أب Country .

هنا نطرح سؤال هل نستخدم ال Snowflake أم ال Star ؟

في ال Star Schema في ال table نفسه نضع ال Cat_id وال Cat_name وال Sc_id وال Sc_name وبهذه الخطوة

ناقصنا مبدأ normalization لأنه أصبح لدينا في كل منتج من هو أباه ومن هو جده وهكذا... بالتالي يوجد هناك

redundant (زيادة) على مستوى ال row .

فيصبح شكل ال product

| Product |
|----------|
| P_id |
| P_name |
| Cat_id |
| Cat_name |
| Sc_id |
| Sc_name |

فمثلاً : لدينا محل يبيع منتجات متنوعة بأقسام متنوعة لمدة 20 سنة وهناك قسم يحوي أدوات منزلية وفيه بعض الأدوات الكهربائية لكنها تندرج تحت Sub_Cat "منزلية" يكون ال rows في product يحمل القيم:

| P_id | P_name | Sc_id | Sc_name | Cat_id | Cat_name |
|------|--------|-------|---------|--------|----------|
| 111 | خلاط | 11 | منزلية | 1 | أدوات |
| 112 | مكسر | 11 | منزلية | 1 | أدوات |

نلاحظ أنه :

في نموذج ERD كانت Category , Sub_Cat كلٍ منها في جدول لوحده (ونفس الأمر بالنسبة لل customer) لذلك لكل منهم record لوحده أما في ال Dimensional يصبحون عبارة عن حقول في جدول product.

قبل ترحيل ال data قام الزبون بإضافة Sub_Cat أبوها "أدوات" واسمها "كهربائية" فيصبح لدينا new record في جدول Sub_Category هو <12 , كهربائية> ثم قام بجعل المنتج "خلاط" ينتمي إلى Sub_Cat "كهربائية" . نفترض أن هذا التغيير طرأ في 2010 فإن مبيعات منتج " خلاط " قبل 2010 مدرج في Sub_Cat "منزلية" وبعد 2010 أصبحت مدرجة ضمن " كهربائية " .

كيف سينعكس هذا الكلام عند نقله الى ال Dimensional في حال القيام بتجميع المنتجات وفق P_id لينتج لدينا total_price وال quantity ونلاحظ أن P_id نفسه في النموذجين؟؟
عندما نقوم بترحيل الداتا سيظهر لدينا ضمن ETL خطأ يان هذا ال record قيمته

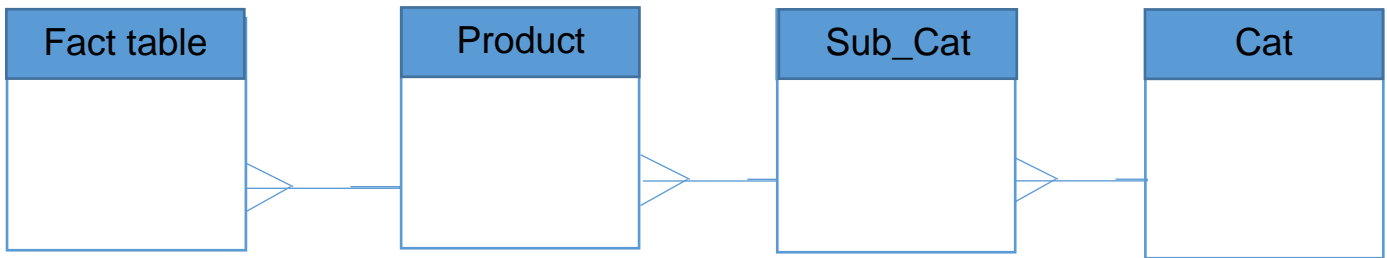
| | | | | | |
|-----|------|----|--------|---|-------|
| 111 | خلاط | 11 | منزلية | 1 | أدوات |
|-----|------|----|--------|---|-------|

سيتضارب مع ال record الجديد

| | | | | | |
|-----|------|----|----------|---|-------|
| 111 | خلاط | 12 | كهربائية | 1 | أدوات |
|-----|------|----|----------|---|-------|

يكون الحل بجعل P_id هو عبارة عن info وتوليد PK عند كل insert في ال Dimensional ونربطه مع Fact_table فإذا أراد معرفة المبيعات من الخلاط ستختلف قبل 2010 وبعدها
ضمنياً نحن نعلم أن كلا ال records السابقين هما لنفس المنتج ولكن كما قلنا سابقاً نحن نقوم بتجميع المنتجات أي اننا سنراكم فوق هذا ال record ، وعندما نصل لمرحلة مشابهة للسابقة أي أننا لا نريد أن نراكم فوق ال record القديم فنضيف عمود يدل ان ال record القديم لم يعد active وأصبح الجديد هو الفعال ويمكن ان أضيف قيمة أخرى تدل على التاريخ الذي أصبح فيه هذا ال record inactive.
ونقوم بهذا الأمر عن طريق ال tools لل ETL في Microsoft و Oracle (وذلك لصعوبة تنفيذ الأمر) وهنا خالفنا مبدأ ال normalization وحافظنا على ال Star Schema .
من الصعب إعادة عملية الترحيل لسنوات سابقة في قاعدة معطيات ضخمة لأنه مكلف زمنياً للترحيل الواحد (والذي نقوم به كل سنة تقريباً) فمن الممكن أن يحتاج ل 3 – 4 أيام .

نمثل الآباء بجداول منفصلة فيكون لدينا عدة مستويات :



الشكل يصبح مثل قطعة ثلج متراكمة فوق بعضها البعض أي انه لم يعد مثل ال Star بويوجد core وهو ال fact table وعدة أبعاد بل أصبح كل بعد لديه ارتباطات أخرى (نفس الشيء بالنسبة لل customer).

ملاحظة :

معظم المستخدمين تميل لل Star Schema من أجل ال performance بالرغم من سلبية وجود redundant في ال data ويعود ذلك إلى أن حجم التخزين الكبير جداً في ال Relational سيتحول لحجم ضئيل في ال Star وفي هذه الحالة لن يؤثر ال Redundant على سير العمل فلن نصل لحجم كبير في ال Star. بينما في حال استخدام Snowflake فهذا يعني وجود عدد من ال Levels مما يجعل الأمر أكثر تعقيد وصعوبة (يوجد الكثير من عمليات join فتكون الداتا موزعة على الجداول).

* في المثال السابق كان لدينا ثلاث records في ال Relational (product, customer, date) وعند الترحيل تجمعت في record واحد في ال Star ضمن ال Fact Table.

في ال MD عند بناء ال DB ننشئ component نسميه hierarchy ونحدد الجدول الذي يعتمد عليه وقد يكون للجدول الواحد n hierarchy فمثلاً في ال date لدينا ال attribute التالية (quarter, year, month, day, week, ...) فممكن أن نبني 1 hierarchy (year, month, day) أو 2 hierarchy (quarter, month) وإمكانية تعدد ال hierarchy هذه احصل عليها كلما زدت ال attribute في ال Dimension.

مثلاً: أريد مبيعاتي من البنطلونات الجينز السوداء ف لدي ال attribute: color.

وإذا احتجت معرفة المبيعات من بنطلونات الجينز ذات القياس M فإنني أقوم بالتفصيل أكثر أي أدخل level أعمق ويمكن ان أصعد 1 level ايضاً باستخدام هذه ال Attribute.

ويمكن ان يتشارك اكثر من cube بنفس ال Dimension وذلك ببناء أكثر من Fact Table.

Written by :

Aisha Awaty

Wordpress and preparation :

Anas Alazmeh

Reviewed by :

Mouayyad Taja

انتهت المحاضرة