Schema Design Notes

Functional Dependencies

A functional dependency is a relationship that exists when **a set of attributes uniquely identifies another set of attributes**. We write that $a_1, a_2, ..., a_n \rightarrow b_1, b_2, ..., b_m$ if when two tuples agree on the attributes $a_1, a_2, ..., a_n$, then they must also agree on the attributes $b_1, b_2, ..., b_m$.

Functional dependencies are crucial for relational schema design. They help eliminate redundancy, which in turn helps reduce **anomalies** such as **redundancy**. They are used to create relations in BCNF form, discussed later.

Some terminology: We can either say a functional dependency **holds** or it does not hold, depending on whether or not it is valid on a specific instance of a relation R. If we can be sure that every instance of R will be one in which the FD holds, we say that R satisfies the FD.

Example: Say we have the following relation, with entry ID (EID) declared a primary key:

EID	Artist	Album	Rating
1	Taylor Swift	Speak Now	10
2	Kanye West	Heartless	2
3	Kanye West	The Life of Pablo	2
4	Charlie Puth	Nine Track Mind	9
5	BTS	Wings	7

Based on the above relation, what valid functional dependencies can we say that we have?

There are actually many possible answers. One valid FD is EID \rightarrow Artist. We can see from the above relation that each matching EID corresponds to the same Artist.

It is also true that EID \rightarrow Album, Rating. We call EID a primary key because it functionally determines all attributes and it is a designated identifier attribute; every EID in the relation is unique. Similarly, we have the valid FDs EID \rightarrow Album, EID \rightarrow Rating, EID \rightarrow Artist, Album, Rating, and so on.

A more interesting example would be Artist \rightarrow Rating. We can look at the table and see that each artist always receives the same rating for their albums. For example, Taylor's albums always receive a 10, Kanye's albums always receive a 2,

Charlie's albums always receive a 9, and Korean-pop sensation BTS' albums always receive a 7. Thus, this functional dependency (Artist \rightarrow Rating) holds.

However, say that another row was added to the relation, consisting of 6 | Taylor Swift | Reputation | 9. If this was the case, the above FD would no longer hold, since one of Taylor's albums now has a rating of 9, and one has a rating of 10 – the ratings don't match!

What would be an example of an FD that does not hold from the get-go? Let's take a look at Artist \rightarrow Album. We can see from the above relation that Kanye West made the albums Heartless and The Life of Pablo. Thus, we have Kanye West \rightarrow Heartless, and Kanye West \rightarrow The Life of Pablo. Remember our definition of valid functional dependencies from before? A given Artist must always correspond to the same Album. However, here we have Kanye corresponding to both Heartless and The Life of Pablo in two separate instances. Thus, we say that this functional dependency does not hold.

As more entries are added to the table, it is possible that some functional dependencies that previously held no longer hold. Likewise, if entries are deleted from the table, some new valid FDs may pop up. For instance, if the row containing The Life of Pablo was deleted, the FD Artist → Album would now hold.

Thus, it is useful to consider FDs that are expected to hold forever, rather than the ones that coincidentally appear in the data.

Closure Algorithm

Given certain functional dependencies, how can we find everything that a set of attributes determine/are a key for? The answer is the closure algorithm. **The closure** algorithm is used for finding all the functional dependencies b1,b2,...,bm for a set of attributes a1,a2,...,an.

```
Algorithm: X = \{A1, ..., An\}. Repeat until X doesn't change do: if (B1, ..., Bn \rightarrow C is a FD) and (B1, ..., Bn are all in X) then add C to X
```

What this is doing is having a **fixed point** and then **tracing the FDs** until there is nothing else to add. This way given a set of attributes we can find all attributes that are in its closure.

The formal notation for specifying a closure is to write $\{a1,a2,...,an\}^{+}$. In the music example, we can write that $\{Artist\}^{+} = \{Artist, Rating\}$.

Keys

When we observe closures, we sometimes get the situation where a closure of a set of attributes is all attributes in a relation. When this happens we can say that the starting set of attributes is called a **superkey**. Of the set of superkeys that are possible for a relation, the superkeys that are the smallest in terms of cardinality are called **minimal keys** or just **keys**. Keys are particularly interesting for database schema design as they indicate what should be marked as primary keys.

Example: We have the same data as before. Let's say that we have the FDs $EID \rightarrow Album$ and $Album \rightarrow Artist$, Rating.

EID	Artist	Album	Rating
1	Taylor Swift	Speak Now	10
2	Kanye West	Heartless	2
3	Kanye West	The Life of Pablo	2
4	Charlie Puth	Nine Track Mind	9
5	BTS	Wings	7

From these given FDs, can we determine what the keys are in the data? One trivial example of a superkey we can get is {EID, Artist, Album, Rating}. But of course this is a superkey! All the attributes in a tuple, will always determine all attributes in the tuple. So generally this is uninteresting.

A more interesting observation is the closure of {EID}. By tracing our FDs, we know that EID determines an Album and an album will determine an artist and rating. Thus, the closure of {EID} is {EID, Artist, Album, Rating} (all attributes in the relation!). Because, we cannot get a smaller size set that is a superkey, {EID} is a minimal key of the relation.

Another minimal key is {Album}. This is a coincidence of the data. Had we added a tuple where different artists recorded an album with the same title, then Album would no longer be a key.

Boyce-Codd Normal Form (BCNF)

Anomalies

When we first gather attributes in a relation to create in a database instance it is possible to not do anything about it and put all your information into a single table. For some data, this is a bad idea due to anomalies. Anomalies are typically through **redundancy**. The typical cause of this is due to there being FDs in a relation that has

the starting attributes not as a key.

Normal Forms and BCNF

Normal forms are a concept in relational models that attempt to promote consistency and ease of use through the partitioning of data into different tables. There are a few different partitioning methods, from the very simple 1st normal form (the only constraint that tables are flat) to the Boyce-Codd Normal Form (BCNF) which is what will be discussed here.

Definition:

A relation is in BCNF if for all non-trivial FDs $X\rightarrow B$, X must be a superkey. Equivalently...

A relation is in BCNF if for all for all X, $\{X\}$ + = X or $\{X\}$ + = {all attributes in the relation}

One can see why BCNF would be a logical choice for how to partition our data as it guarantees that there are no anomalies as discussed before.

There are systematic algorithms to convert non-normalized data to normal forms like BCNF given any functional dependencies that you define. For BCNF in particular, lossless decomposition is the methodology for partitioning, and the chase algorithm is used for verification, however, you are not expected to know how these algorithms work for the final.

Takeaway

What we hope that you take away from this small introduction to design theory for databases is that there are concrete ways to make a well-behaved schema. In any future applications that you may make which use relational databases, if you are able to correctly define your functional dependencies (properties of your data), you can normalize your model to promote consistency and ease of use.

Material adapted from CSE 344 AU 2017 slides

Written by: Zealous TAs
Allison Chou (aachou @ cs)
Jonathan Leang (jleang @ cs)
Shana Hutchison (shutchis @ cs)