Washington State University CPT_S 415 – Big Data Online

Srinivasulu Badri

Assignment 1

Name: Nam Jun Lee

Student Number: 11606459

1. [**Big Data concept**] Give one example of Big Data application you know. Use the detailed example to explain each of the five Big V's.

The big data application I know is Facebook. So let me explain how 5V's of big dat a are applied on Facebook.

Five Big V's:

- (1) Volume: Facebook has a huge number of users. According to Google, monthly active users are 2.934 billion. As such, the amount of data Facebook has is enormous.
- (2) Velocity: Facebook speeds up scrolling posts and playing videos, so users don't feel uncomfortable using the service
- (3) Variety: Facebook contains text such as users' personal information and messag es. There are also pictures and countless videos. As such, Facebook includes v arious data source types.
- (4) Veracity: Facebook contains the quality and accuracy of its data. They remove dirty data for the smooth operation of the service so that users do not feel un pleasant emotions while using the service.
- (5) Value: Facebook includes a 'Like' feature. Through the information obtained by this function, the product advertisement for each user's area of interest is disp layed intensively, thereby arousing the desire of users to purchase. As such, it helps to create business value by analyzing users' interests.
- 2. [Relational Data Model] As of January 2017, the OpenFlights Airports Databa se (https://openflights.org/data.html) contains over 10,000 airports, train stations and ferry terminals spanning the globe. Each entry in the Airport table contain s the following:
 - A. Consider the following terms: relation schema, relational database sche ma, domain, attribute, attribute domain, relation instance. Give what the se terms are with the above Airport database. Give one small (4-5 tuple s) instance of the Airport table.

Airport

Airport	Name	City	Countr	IATA	ICAO	Latitude	Longit	Altit	Tim	DST	Tz	Type	Source
ID			у				ude	ude	e		database		
									zone		time		
											zone		
1	Goroka	Goroka	Papua	GKA	AYGA	-6.0816	145.39	5282	10	U	Pacific/P	airport	OurAirports
	Airport		New			8983459	19982				ort_Mor		
			Guinea			0001	91				esby		
2	Madang	Madang	Papua	MAG	AYM	-5.2070	145.78	20	10	U	Pacific/P	airport	OurAirports
	Airport		New		D	7988739	90014				ort_Mor		
			Guinea				65				esby		
3	Mount	Mount	Papua	HGU	AYM	-5.826	144.29	5388	10	U	Pacific/P	airport	OurAirports
	Hagen	Hagen	New		Н	7898559	60052				ort_Mor		
	Kagamuga		Guinea			57031	49023				esby		
	Airport						44						
4	Nadzab	Nadzab	Papua	LAE	AYNZ	-6.569	146.72	239	10	U	Pacific/P	airport	OurAirports
	Airport		New			803	5977				ort_Mor		
			Guinea								esby		

Relation schema: Airport(AirportID, Name, City, Country, IATA, ICAO, Latitude, Longitude, Altitude, Time zone, Tz database time zone, Type, Source)

Relational database schema: Database schema is the collection of relation schema. Hence, one relational database schema.

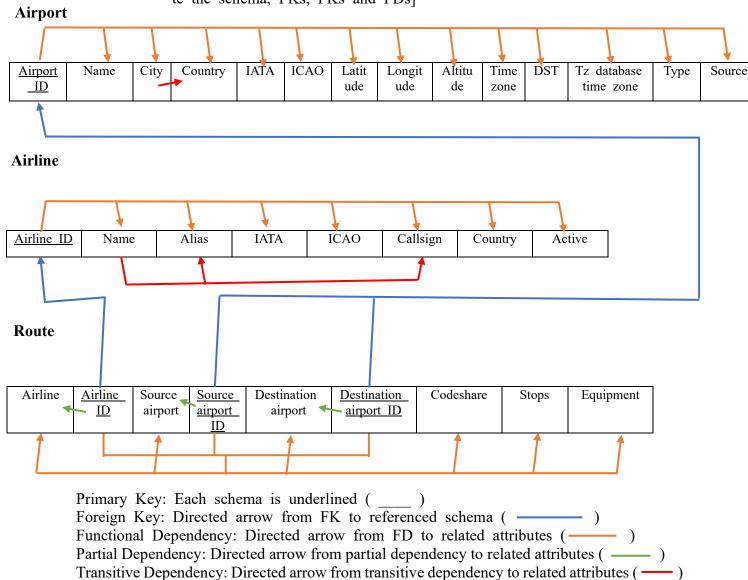
Domain: A data type with optional constraints. For example, the domain of Name, City, and Country in the Airport relation should be a string. It cannot have numeric values.

Attribute: Column Headers. Hence, Airport table has 14 attributes: AirportID, Name, City, Country, IATA, ICAO, Latitude, Longitude, Altitude, Time zone, Tz database time zone, Type, Source.

Attribute domain: All possible values in a Column. For example, the domain of IATA in the Airport relation is 3-letter IATA code.

Relation instance: A relation instance is a tuple/row in a relation. Hence, the Airport table has one particular combination of attribute values.

B. There are three databases in the OpenFlight dataset: Airport, Airline, an d Route. Give the schema of these three databases and mark the primar y keys, foreign keys and provide examples of functional dependencies y ou identified over the three tables. [You may draw a diagram to illustra te the schema, PKs, FKs and FDs]



- 3. [Functional Dependencies] Recall Armstrong's axioms.
- Reflexivity rule: if $Y \subseteq X$ then $X \to Y$
- Augmentation rule: if $X \rightarrow Y$ then $XZ \rightarrow YZ$
- Transitivity rule: if $X \to Y$ and $Y \to Z$ then $X \to Z$
 - a. Give two examples for using Armstrong's inference rules to induce new FDs from the set of FDs you designed in question 2 (b).
 - (1) Augmentation rule: Airport Schema

If Airport ID
$$\rightarrow$$
 Country then,

$$Airport ID + Name = Country + Name$$

(2) Transitivity rule: Airline Schema

If Alias
$$\rightarrow$$
 Name and Name \rightarrow Callsign then,

- b. Prove the following inference rules also hold, using FD definition and Armstrong's Axioms.
 - i. decomposition rule: if $X \to YZ$ then: $X \to Y$ and $X \to Z$

Proof:

$$X \rightarrow YZ$$
 ---- given above

$$YZ \rightarrow Y$$
 ----- Reflexivity Rule

$$YZ \rightarrow Z$$
 ----- Reflexivity Rule

$$X \rightarrow Y$$
 ----- Transitivity Rule on $X \rightarrow YZ$ and $YZ \rightarrow Y$

$$X \rightarrow Z$$
 ---- Transitivity Rule on $X \rightarrow YZ$ and $YZ \rightarrow Z$

Hence, Proved.

ii. Pseudo transitivity: if $X \rightarrow Y$ and $YW \rightarrow Z$ then: $XW \rightarrow Z$

Proof:

$$X \rightarrow Y$$
 ---- given above

$$YW \rightarrow Z$$
 ---- given above

$$XW \rightarrow YW$$
 ----- Augmentation Rule on $X \rightarrow Y$ by augmenting with W

 $XW \rightarrow Z$ ----- Transitivity Rule on $XW \rightarrow YW$ and $YW \rightarrow Z$

Hence, Proved.

4. [Normalization] Given a relation R(A₁, A₂, A₃, A₄), with three FDs A₂, A₃ \rightarrow A₄; A₃, A₄ \rightarrow A₁; A₁, A₂ \rightarrow A₃. Provide the 3NF and BCNF form of the schema and explain why.

First, find the candidate keys.

Check FDs: A_1 , A_3 , A_4 are included in both attributes, but A_2 is included only in the left attribute. So, calculate the closing of the power set combination of A_2 :

$$A_1, A_2 = A_1, A_2, A_3, A_4 = R$$
 ---- Candidate key

$$A_2, A_3 = A_1, A_2, A_3, A_4 = R$$
 ---- Candidate key

$$A_2, A_4 = A_2, A_4 \subset R$$

 $\{A_1, A_2\}, \{A_2, A_3\}$ lead only to superkeys, so $\{A_1, A_2\}, \{A_2, A_3\}$ are candidate keys.

Hence, the Set of key attributes: A_1 , A_2 , A_3 .

Now, provide the 3NF and BCNF form of this relation R:

- 3NF: For each FD, determine whether the left-hand side is a super key or the right-hand side is all key attributes, **this relation is in 3NF** because there is no transitive dependency here.
- BCNF: Check each FD: $\{A_3, A_4\}$ is non-trivial but is not a super key. If this relation is in BCNF, each non-trivial FD must be a super key. Hence, **it violates BCNF**. It means **this relation is not in BCNF**.