

Data Insight

Chapter 1: Data Cleaning

- **Descriere**

Scopul proiectului nostru este de a permite utilizatorilor să realizeze și vizualizeze graficele aferente statisticilor efectuate pe seturi mari de date.

Inițial intenționăm să realizăm o multitudine de funcții care vor extrage date dintr-o bază de date oarecare și vor crea diferite statistici pe baza acestor date.

Ulterior am realizat că acest lucru nu este foarte greu de realizat, noi trebuind doar să implementăm în `plsql` anumite funcții deja existente în cărți/pe internet și să le facem să se poată adapta pe mai multe tipuri de tabele. Pe lângă asta, în cazul în care datele dintr-un tabel n-ar fi perfect corect inserate, funcțiile noastre ar genera erori.

Astfel, scopul real al proiectului nostru a devenit realizarea primului pas ce trebuie făcut în momentul în care dorim să generăm statistici pentru o anumită bază de date, și anume **curățarea datelor din tabele**.

Curățarea datelor unei tabele implică numeroase verificări care, în mod normal, sunt extrem de dificil de făcut chiar și într-o companie foarte mare. Acest proces durează extrem de mult, iar bazele de date nu vor fi, în general, curățate complet. De asemenea, curățarea manuală a datelor din baza de date a unei aplicații aflată în producție este aproape imposibil de realizat, din cauza multitudinii de date. În principiu, se poate spune că se poate evita inserare greșită a datelor prin realizarea de verificări automate în momentul inserției, dar acoperirea tuturor cazurilor ce implică date greșite este, aproape imposibil de realizat, chiar și în mod automat.

Aplicația noastră va avea următoarele funcții de curățare a datelor, funcții ce se vor aplica pe orice structură de tabele:

- ❖ **Rezolvarea duplicatelor:**

Existența duplicatelor este una din cele mai întâlnite probleme în bazele de date. De exemplu pot exista două (sau mai multe) înregistrări ale aceleiași persoane într-un tabel cu același nume, prenume și adresă de email, dar cu numere de telefon diferite pentru că validarea la înregistrare a fost făcută după numărul de telefon. Pentru aplicațiile unde nu este permisă utilizarea a mai mult de un cont per user acest fenomen ar fi o adevărată problemă. Ca și soluții vom aborda: merge între înregistrări, ștergerea duplicatelor.

- ❖ **Conversia datelor la tipul corespunzător:**

Într-o coloană a unui tabel putem avea numai înregistrări de tip `number`, dar acea coloană să aibă ca tip `varchar`. Astfel se va face recunoașterea tipurilor de date înregistrate în tabel și validarea acestora cu tipul coloanelor, urmând modificarea corespunzătoare a tipului coloanelor în cazul în care este necesar. Totodată dacă avem o coloană de tip `varchar(100)`, dar noi avem numai înregistrări de maxim 10 caractere în acea coloană vom reduce dimensiunea acelui `varchar`.

❖ **Rezolvarea erorilor de sintaxă:**

Erorile de sintaxă sunt de mai multe tipuri. Vom încerca să le rezolvăm pe următoarele:

- white spaces: spațiile albe de la începutul și sfârșitul înregistrărilor trebuie șterse (trim)
- fix typos: în anumite situații ne dorim să facem selecția datelor după anumite standarde, dar datele pot avea mici deviații de la standardele propuse; un exemplu tipic este situația în care dorim să extragem userii având sexul bărbat dintr-o tabelă folosind tagul "Male", dar noi avem bărbații înregistrați în tabelă sub următoarele taguri: Man, male, MALE, M., masculin, bărbat.
- pad strings: uneori ne dorim ca anumite tipuri de date să aibă o dimensiune fixă, putem completa astfel înregistrările cu spații albe la final în cazul stringurilor sau cu 0-uri la început în cazul numerelor până când ajungem la dimensiunea dorită

❖ **Standardizarea datelor:**

Într-o tabelă ne dorim să avem datele introduse după anumite standarde:

- numele și prenumele încep cu majusculă, iar restul caracterelor sunt minuscule (Andrei, NU aNDREI/andrei/ANDREI)
- nu avem majuscule intercalate cu minuscule (mașină/MAȘINĂ, nu mAșiNĂ)
- vrem să avem data de același format peste tot (zz.ll.aaaa, NU ll/zz/aa sau aaaa-ll-zz)

❖ **Scalarea/Transformarea datelor:**

Vrem să oferim posibilitatea de a transforma toate datele sau datele 'greșite' în formatul dorit. De exemplu, într-un tabel noi avem notele studenților înregistrate, dar aceste note sunt puse în două sisteme de notare diferite: un sistem este cu note de la 1 la 5, iar celălalt de la A la F. Vom notele astfel încât acestea să fie într-un singur sistem de notare.

❖ **Normalizarea datelor:**

Normalizarea datelor este un proces ce ajută extrem de mult la generarea statisticilor. Pentru realizarea anumitor statistici este nevoie ca datele să fie distribuite uniform și la o scală cât mai mică posibil pentru eficientizarea funcțiilor generatoare de statistici.

Totuși, crearea unei aplicații care să curețe în mod automat și complet întreagă bază de date necesită cunoștințe mult prea avansate și un volum extrem de mare de muncă. Așadar aplicația noastră va trebui puțin configurată înainte de execuția algoritmilor de curățare:

- Se va selecta operația de curățare care se dorește a fi executată
- Se va specifica ce date vor fi modificate și în ce mod (în cazul standardizării, de exemplu)
- Se va specifica dacă se dorește doar vizualizarea erorilor existente sau rezolvarea directă a lor cu o soluție prestabilită/setată de utilizator.

-Va exista și o serie de funcții predefinite care vor recunoaște după anumite key tipurile coloanelor, executând automat curățarea completă a datelor din acele coloane
Exemplu:

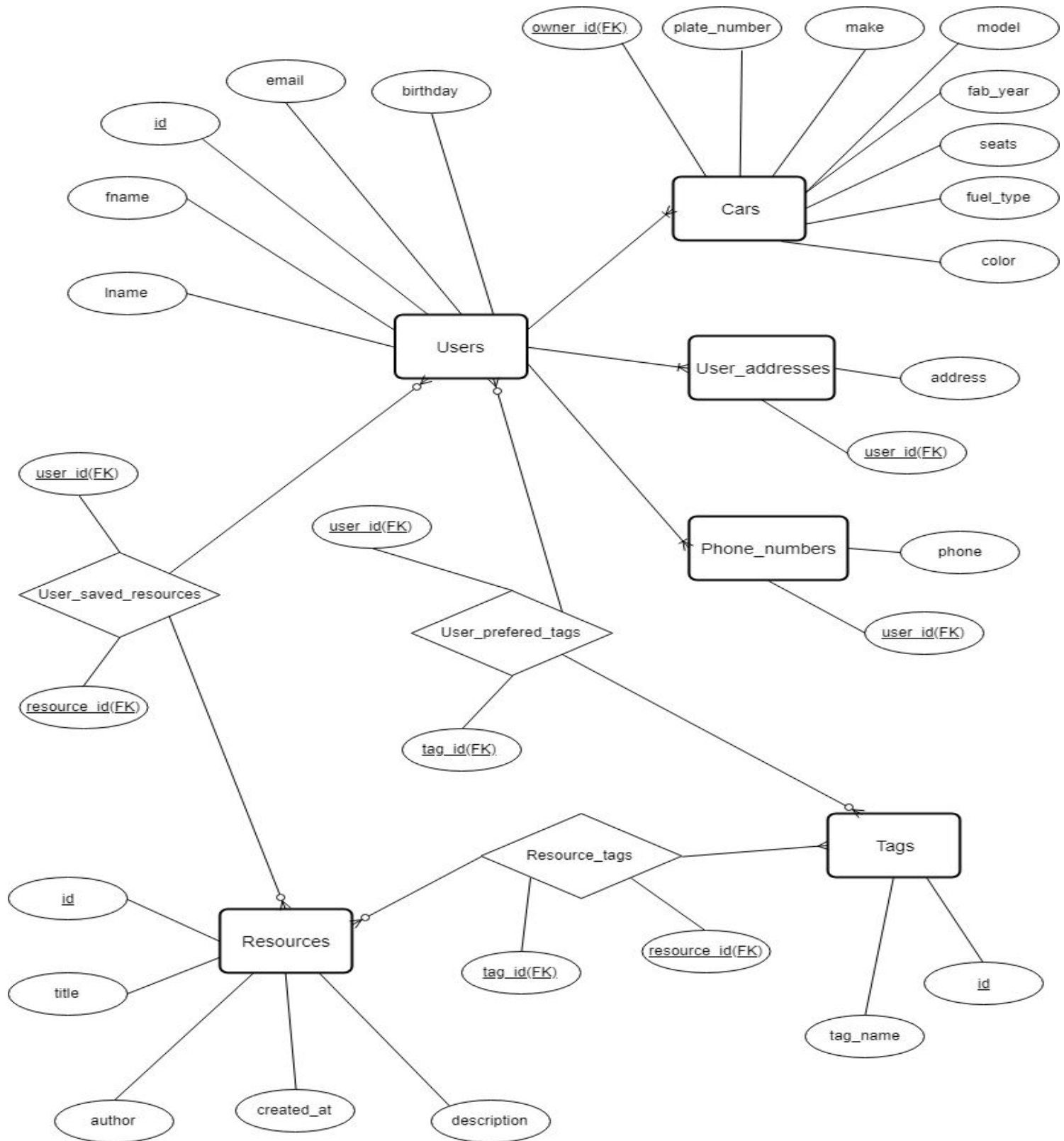
- standardizarea numerelor de telefon în formatul +[prefix țară]număr, prefixul de țară fiind automat realizat în funcție de o coloană în care a fost introdusă țara;

- standardizarea adresei în formatul City [city], Str. [str], No. [nr]

- semnalarea erorilor în cazul unei coloane cu note, în care apar note mai mari decât 10.

Pentru testarea aplicației am ales să creăm o bază de date cu o tabelă ce are diverse tipuri de date și să generăm înregistrări cu anumite tipuri de erori pentru a le putea curăța utilizând funcționalitatea algoritmilor noștri.

- Schemă Entitate/Relație



- **Formă normală + Argumentare**

Schemele relaționale și dependențele funcționale

Pornind de la tabelele considerate pentru a modela aplicația, vom încerca să ducem baza de date în forma normală 4. Așadar, inițial avem relațiile:

USERS (id, email, fname, lname, birthday, phone, address, car_plate_number, car_make, car_model, car_fab_year, car_seats, car_fuel_type, car_color, preferred_tags)

Dependențe funcționale USERS:

$id \rightarrow X, X \in \text{USER}$

$car_plate_number \rightarrow Y, Y \in \{car_make, car_model, car_fab_year, car_seats, car_fuel_type, car_color\}$

RESOURCES (id, author, title, created_at, description, tags)

Dependențe funcționale RESOURCES:

$id \rightarrow X, X \in \text{RESOURCES}$

1NF

Cunoscând faptul că pe viitor aplicația ar putea oferi suport pentru multiple numere de telefon, ceea ce ar încălca prima formă normală, extragem câmpul într-o tabelă separată. Obținem astfel:

USERS (id, email, fname, lname, birthday, address, car_plate_number, car_make, car_model, car_fab_year, car_seats, car_fuel_type, car_color, preferred_tags)
PHONE_NUMBERS (user_id (fk), phone)

Același lucru se aplică și pentru câmpul preferred_tags, acesta fiind plural. Obținem astfel:

USERS (id, email, fname, lname, birthday, address, car_plate_number, car_make, car_model, car_fab_year, car_seats, car_fuel_type, car_color)
PHONE_NUMBERS (user_id (fk), phone)
USER_PREFERRED_TAGS (user_id (fk), tag)

În mod analog punctului unu, preconizăm faptul că un utilizator ar putea avea multiple adrese, iar în cazul în care aplicația ar avea nevoie să cunoască adresa efectivă a unui utilizator, se poate crea o tabelă de asociere, cu cheie străină către adresa considerată efectivă. Considerăm faptul că adresa este atomică, iar aceasta va fi preluată dintr-un singur câmp. După extragerea câmpului adresă din relația utilizatori, obținem:

USERS (id, email, fname, lname, birthday, address, car_plate_number, car_make, car_model, car_fab_year, car_seats, car_fuel_type, car_color)

$id \rightarrow X, X \in \text{USER}$

$car_plate_number \rightarrow Y, Y \in \{car_make, car_model, car_fab_year, car_seats, car_fuel_type, car_color\}$

PHONE_NUMBERS (user_id (fk), phone)

$id \rightarrow phone$

USER_PREFERED_TAGS (user_id (fk), tag_name)

$id \rightarrow tag$

USER_ADDRESSES (user_id (fk), address)

$id \rightarrow address$

Analog punctului doi, resursele vor avea multiple taguri, deci trebuie sa extragem din relație câmpul tags și să creem o altă relație numită resources tags.

RESOURCES (id, author, title, created_at, description)

Dependențe funcționale RESOURCES:

$id \rightarrow X, X \in \text{RESOURCES}$

RESOURCE_TAGS (resource_id (fk), tag_name)

Observăm faptul că avem doua relatii similare: USER_PREFERED_TAGS, respectiv RESOURCE_TAGS, având un câmp comun, tag_name. Dacă în viitor, cineva ar dori să vadă ce taguri sunt în sistem, ar trebui să interogheze doua tabele. Așadar putem crea o noua relatie, numita Tags, la care vor face referire cheile straine din cele doua tabele. Un alt avataj al acestei acțiuni este reprezentat de faptul că dacă un utilizator preferă un tag ce aparține in mod întâmplător unei resurse, denumirea acestui tag va fi consistentă atat la utilizator, cât și in resursa. Continuând cu modificările de relatii obținem:

USERS (id, email, fname, lname, birthday, car_plate_number, car_make, car_model, car_fab_year, car_seats, car_fuel_type, car_color)

$id \rightarrow X, X \in \text{USER}$

$car_plate_number \rightarrow Y, Y \in \{car_make, car_model, car_fab_year, car_seats, car_fuel_type, car_color\}$

PHONE_NUMBERS (user_id (fk), phone)

USER_PREFERED_TAGS (user_id (fk), tag_id (fk))

USER_ADDRESSES (user_id (fk), address)

RESOURCES (id, author, title, created_at, description)

$id \rightarrow X, X \in \text{RESOURCES}$

RESOURCE_TAGS (resource_id (fk), tag_id (fk))

TAGS (id, tag_name)

2NF

Notăm cu $k(R)$ mulțimea cheilor candidat a relației R.

$k(\text{USERS}) = \{(id, car_plate_number)\}$

$k(\text{RESOURCES}) = \{id\}$

$k(\text{TAGS}) = \{\text{id}\}$

Observăm că în relația USERS exista o cheie candidat, iar din dependențele funcționale rezultă faptul că attributele car_make, car_model, car_fab_year, car_seats, car_fuel_type, car_color sunt dependente parțial de submulțimea cheii candidat, formată doar din car_plate_number. Așadar, trebuie să extragem toate attributele ne-prime în altă tabelă și să referențiem user-ul. După procedeu, obținem:

USERS (id, email, fname, lname, birthday)

$\text{id} \rightarrow X, X \in \text{USER}$

CARS (owner_id (fk), plate_number, make, model, fab_year, seats, fuel_type, color)

$\text{car_plate_number} \rightarrow Y, Y \in \{\text{car_make}, \text{car_model}, \text{car_fab_year}, \text{car_seats}, \text{car_fuel_type}, \text{car_color}\}$

Recalculând mulțimile de chei candidat a relației R, obținem:

$k(\text{USERS}) = \{\text{id}\}$

$k(\text{CARS}) = \{\text{plate_number}\}$

$k(\text{RESOURCES}) = \{\text{id}\}$

$k(\text{TAGS}) = \{\text{id}\}$

Deoarece mulțimile de chei a relației R au fiecare doar un singur element, putem să ne oprim din evaluat și să concluzionăm că schema de relație R este în 2NF.

3NF

Reevaluăm schema de relație R, încercând să verificăm dacă toate attributele ne-prime sunt ne-transitiv dependente de orice cheie candidat a unei relații.

USERS (id, email, fname, lname, birthday)

$\text{id} \rightarrow X, X \in \text{USER}$

CARS (owner_id (fk), plate_number, make, model, fab_year, seats, fuel_type, color)

$\text{car_plate_number} \rightarrow Y, Y \in \{\text{car_make}, \text{car_model}, \text{car_fab_year}, \text{car_seats}, \text{car_fuel_type}, \text{car_color}\}$

PHONE_NUMBERS (user_id (fk), phone)

USER_PREFERRED_TAGS (user_id (fk), tag_id (fk))

USER_ADDRESSES (user_id (fk), address)

RESOURCES (id, author, title, created_at, description)

$\text{id} \rightarrow X, X \in \text{RESOURCES}$

RESOURCE_TAGS (resource_id (fk), tag_id (fk))

TAGS (id, tag_name)

Deoarece toate attributele ne-prime din relațiile de mai sus depind funcțional direct de cheie și nu și de un alt atribut ne-prim, rezultă că schema de relație R este în 3NF.

BCNF

Calculând mulțimile de chei candidat a relației R, obținem:

$k(\text{USERS}) = \{\text{id}\}$

$k(\text{CARS}) = \{\text{plate_number}\}$

$k(\text{RESOURCES}) = \{\text{id}\}$

$k(\text{TAGS}) = \{\text{id}\}$

Datorită faptului că pentru orice dependență funcțională $X \rightarrow Y$ din schemele de relație USERS, CARS, RESOURCES, PHONE_NUMBERS, USER_PREFERRED_TAGS, USER_ADDRESSES, RESOURCE_TAGS, TAGS, X este supercheia schemei R sau $X \rightarrow Y$ este o dependență funcțională trivială, rezultă că schema de relație R este în BCNF.

4NF

Schema de relație R și dependențele multivaluate:

USERS (id, email, fname, lname, birthday)

$\text{id} \twoheadrightarrow X, X \in \text{USER}$

CARS (owner_id (fk), plate_number, make, model, fab_year, seats, fuel_type, color)

$\text{car_plate_number} \twoheadrightarrow Y, Y \in \{\text{car_make, car_model, car_fab_year, car_seats,}$

$\text{car_fuel_type, car_color}\}$

PHONE_NUMBERS (user_id (fk), phone)

USER_PREFERRED_TAGS (user_id (fk), tag_id (fk))

USER_ADDRESSES (user_id (fk), address)

RESOURCES (id, author, title, created_at, description)

$\text{id} \twoheadrightarrow X, X \in \text{RESOURCES}$

RESOURCE_TAGS (resource_id (fk), tag_id (fk))

TAGS (id, tag_name)

Calculând mulțimile de chei candidat a relației R, obținem:

$k(\text{USERS}) = \{\text{id}\}$

$k(\text{CARS}) = \{\text{plate_number}\}$

$k(\text{RESOURCES}) = \{\text{id}\}$

$k(\text{TAGS}) = \{\text{id}\}$

Datorită faptului că pentru orice dependență multivaluată $X \twoheadrightarrow Y$ din schemele de relație USERS, CARS, RESOURCES, PHONE_NUMBERS, USER_PREFERRED_TAGS, USER_ADDRESSES, RESOURCE_TAGS, TAGS, X este supercheia schemei R sau $X \twoheadrightarrow Y$. Așadar, putem spune că schema de relație este în 4NF.

- **Script creare db + populare**
- **Descriere și motivare structuri folosite**

Deoarece aplicația noastră trebuie să funcționeze pe diverse tipuri de baze de date și, implicit, tabele este foarte dificil să folosim indecși sau views și nu ar avea nici o utilitate prea mare din moment ce parcurgerea datelor dintr-un tabel se va face doar o dată, sau, cel mult de două ori, în cazul în care utilizatorul dorește mai întâi să vadă erorile și apoi să curețe tot/să aleagă ce să curețe. Pe lângă asta nu avem nevoie de eficiență prea ridicată

din moment ce aplicația va servi câte unui singur client (nu este server/app web, ci doar un tool) care își va alege funcțiile de curățare pe care va dori să le ruleze pe baza lui de date.

Pentru testarea aplicației vom folosi tabelele dintr-o aplicație uzuală (aplicația de la Tehnologii Web care aduce articole de pe mai multe site-uri IT related și le afișează în funcție de cerințele utilizatorului), la care am mai adăugat o tabelă Cars pentru a varia tipurile de date pentru testare. Ulterior vom mai adăuga diferite tabele/baze de date existente pe internet pentru a verifica corectitudinea algoritmilor de curățare utilizați.

Tabele folosite inițial:

USERS: va reține datele utilizatorilor.

USER_ADDRESSES: va reține adresele utilizatorilor (un utilizator poate avea una sau mai multe adrese)

PHONE_NUMBERS: va reține numerele de telefon ale utilizatorilor (un utilizator poate avea unul sau mai multe numere de telefon)

RESOURCES: aici vor fi stocate temporar informațiile despre articolele aduse din diverse surse

TAGS: aici vor fi reținute tag-urile disponibile în aplicație

USER_PREFERRED_TAGS: aici vor fi stocate tag-urile preferate de utilizatori. Un utilizator poate avea 0 sau mai multe tag-uri preferate, iar un tag poate fi atribuit pentru 0 sau mai mulți utilizatori.

USER_SAVED_RESOURCES: această tabelă va reține resursele adăugate la salvate pentru fiecare utilizator. Un utilizator poate avea 0 sau mai multe resurse salvate, iar o resursă poate fi atribuită pentru 0 sau mai mulți utilizatori.

RESOURCES_TAGS: aici vor fi stocate tag-urile atribuite fiecărei resurse. O resursă poate avea 1 sau mai multe tag-uri atribuite, iar un tag poate fi atribuit pentru 0 sau mai multe resurse.

CARS: aceasta este tabela pe care am adăugat-o pentru a varia tipurile de date pe care să testăm aplicația noastră. tabela va fi legată de utilizatori, aceștia având posibilitatea de a avea o mașină și va reține toate informațiile despre mașină.

● Listă funcții + descriere

→ *mergeDuplicates*: Funcția va rezolva duplicatele din tabelele noastre făcând merge între ele după anumite criterii (unul dintre ele ar fi relația acestor date cu datele din alte tabele) -> se vor uni duplicatele într-o singură înregistrare, iar datele din alte tabele care aveau ca foreign key unul din câmpurile înregistrărilor care au făcut merge își vor schimba foreign key-ul (dacă e necesar).

→ *deleteDuplicates*: Funcția va păstra una din înregistrări, în funcție de anumite criterii pe care le va selecta userul - cele mai puține câmpuri null, cele mai multe aparențe ale unui câmp ca foreign key în alte tabele, random - și va șterge toate celelalte duplicate.

→ *resolveTypeDimension*: Funcția va modifica dimensiunile declarate ale tipurilor de date ale coloanelor astfel încât să nu fie alocate dimensiuni prea mari. -> de exemplu pentru cazul în care avem declarat varchar(100), dar noi avem numai înregistrări de maxim 10 caractere, dimensiunea se va modifica la 11, sau la 15-20 dacă ne dorim și o marjă de eroare

→ *resolveTypeConflicts*: Funcția va reface tipurile de date declarate coloanelor

astfel încât să se plieze strict pe înregistrările noastre (de exemplu dacă noi avem înregistrate numai numere, dar tipul coloanei este varchar, tipul coloanei va deveni number)

→ *deleteUnnecessaryWhiteSpaces*: Funcția va șterge spațiile în plus de la începutul/sfârșitul înregistrărilor

→ *fixTypos*: Funcția va repara stringurile care ar trebui să fie uniformizate. De exemplu în loc de stringul "Male" avem introduse "Masc., M., male". Funcția va primi ca parametru coloana și condițiile de uniformizare.

→ *addStringsPadding*: Funcția va completa înregistrările pentru care ne dorim o dimensiune fixă. Pentru stringuri va adăuga spații la final, iar pentru numere va adăuga 0-uri la început.

→ *dataStandardize*: Funcția va standardiza datele după anumite standarde primite ca parametru. De exemplu va standardiza adresa după formatul City [city] Str. [str] No. [nr]

→ *dataTransform*: Funcția va transforma datele (toate, sau cele care specificate) dintr-un format în altul.

→ *dataNormalize*: Funcția va normaliza datele din tabelele/câmpurile specificate, astfel încât să fie distribuite uniform în tabel, sau să fie scalate la limite mai mici.

→ *autoCleaning*: Funcția va curăța automat toată baza de date selectată după anumite criterii setate default din funcție. Aceasta va face recunoașterea anumitor tipuri de câmpuri (Name, Phone, Address etc) și a tipurilor și va încerca să curețe cât de mult posibil baza de date.