

# “Requêter une base de données avec SQL”

Ce projet a pour objectif de créer et gérer une base de données pour analyser le marché des assurances habitation en France avec SQL.

## Objectifs

- Comprendre les types de données utilisées.
- Créer un schéma relationnel pour structurer les données.
- Création d'une base de données.
- Charger les données dans un SGBD et vérifier leur intégrité.
- Réaliser des analyses avec des requêtes SQL.

# 1.1 Exploration des données

Dictionnaire des données: "CONTRAT.CSV"

## Données sur les contrats d'assurance habitation

Nom des colonnes	Type de données	Taille	Clé	Description
<b>Contrat_ID</b>	VARCHAR	8	PK	ID unique pour les contrats. Les données sont seulement des chiffres avec une taille de 6.
<b>No_voie</b>	VARCHAR	5		Numéro de voie pour l'adresse du logement assuré.
<b>B_T_Q</b>	VARCHAR	1		Indicateur éventuel de répétition pour l'adresse du logement assuré sur un caractère
<b>Type_de_voie</b>	VARCHAR	5		Abbréviation du type de voie pour l'adresse du logement assuré: RUE, AV (Avenue), QUAI, ecc
<b>Voie</b>	VARCHAR	40		Libellé de la voie pour l'adresse du logement assuré
<b>Code_dep_code_commune</b>	VARCHAR	6	PFK	Concaténation du code départemental officiel avec le code commune pour avoir une clé unique
<b>Code_postal</b>	CHAR	5		Code postal pour l'adresse du logement assuré

Nom des colonnes	Type de données	Taille	Clé	Description
<b>Surface</b>	MEDIUMINT UNSIGNED	5		La surface du logement assuré en mètres carrés
<b>Type_local</b>	VARCHAR	20		Le type du bien assuré (appartement, maison.. Ecc)
<b>Occupation</b>	VARCHAR	35		La personne qui occupe le bien (locataire, propriétaire... etc.)
<b>Type_contrat</b>	VARCHAR	50		Le type de contrat que l'assurance a stipulé avec le client.
<b>Formule</b>	VARCHAR	15		La formule choisie par le client avec le contrat stipulé
<b>Valeur_declaree_biens</b>	VARCHAR	12		La valeur déclarée pour les biens.
<b>Prix_cotisation_mensuel</b>	MEDIUMINT UNSIGNED	5		Le prix de la cotisation mensuelle pour le client

# 1.2 Exploration des données

Dictionnaire des données: "REGION.CSV"

Référentiel des régions françaises extrait de *data.gouv.fr*

Nom des colonnes	Type de données	Taille	Clé	Description
<b>Code_dep_code_commune</b>	VARCHAR	6	PK	Concaténation du code départemental officiel avec le code commune pour avoir une clé unique.
<b>reg_code</b>	VARCHAR	2		L'identifiant régional officiel sans préfixe.
<b>reg_nom</b>	VARCHAR	26		Le nom officiel de la région.
<b>aca_nom</b>	VARCHAR	30		Le centre académique de la zone géographique.
<b>dep_nom</b>	VARCHAR	43		Le nom officiel du département.
<b>com_nom_maj_court</b>	VARCHAR	45		Le nom de la ville.
<b>dep_code</b>	VARCHAR	3		Le code départemental officiel
<b>dep_nom_num</b>	VARCHAR	28		<i>Le nom du département avec le code départemental</i> (redondance: c'est la concatenation de "dep_nom" + "dep_code")

Chaque type de donnée a été associé à une contrainte de taille spécifique afin d'optimiser l'efficacité de la base de données et d'accélérer les requêtes et manipulations.

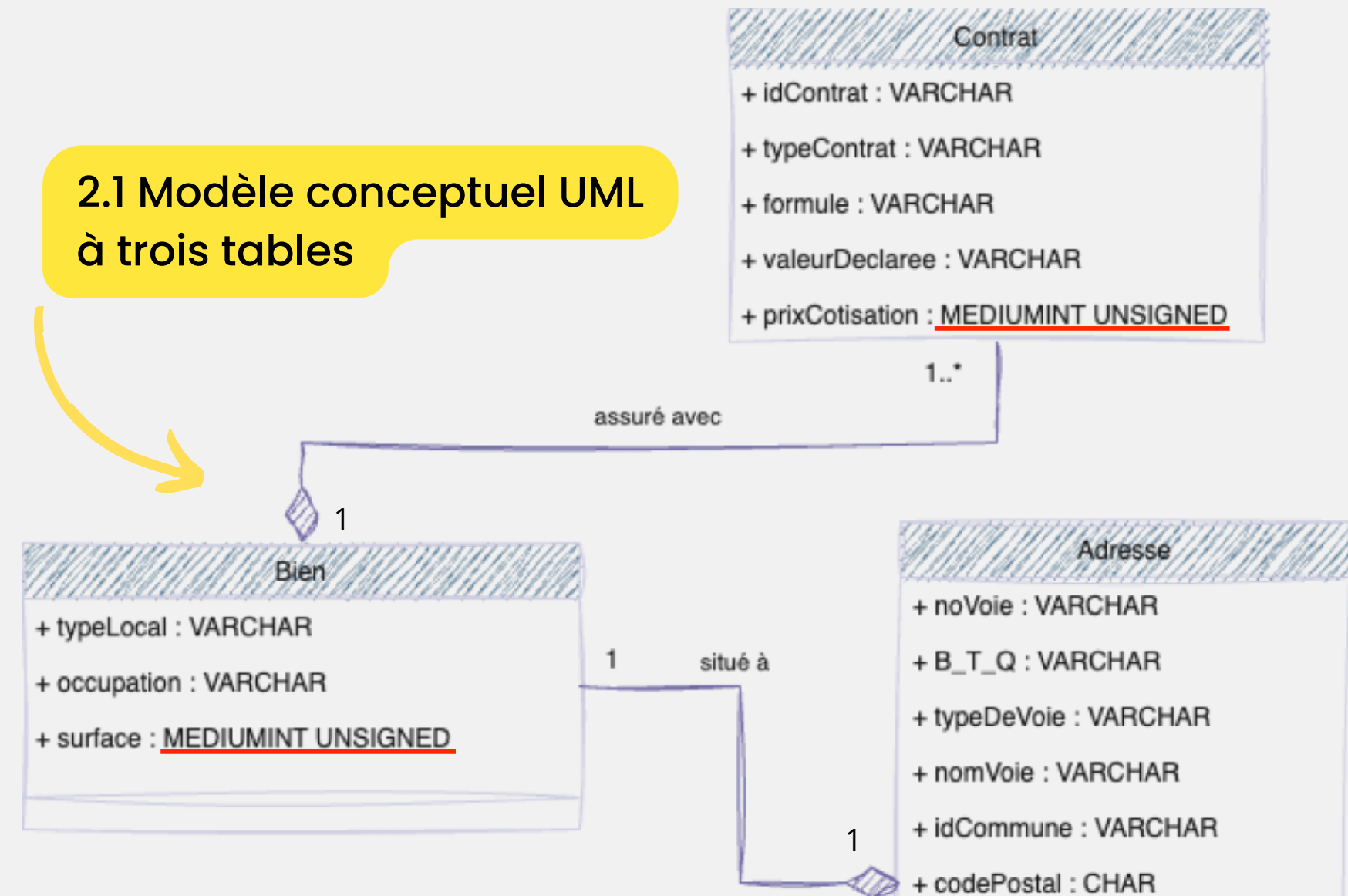
Cette démarche s'appuie sur une recherche approfondie, prenant en compte toutes les possibilités d'insertion et une analyse minutieuse des données existantes.

# 2. Modèle conceptuel UML et modèle relationnel

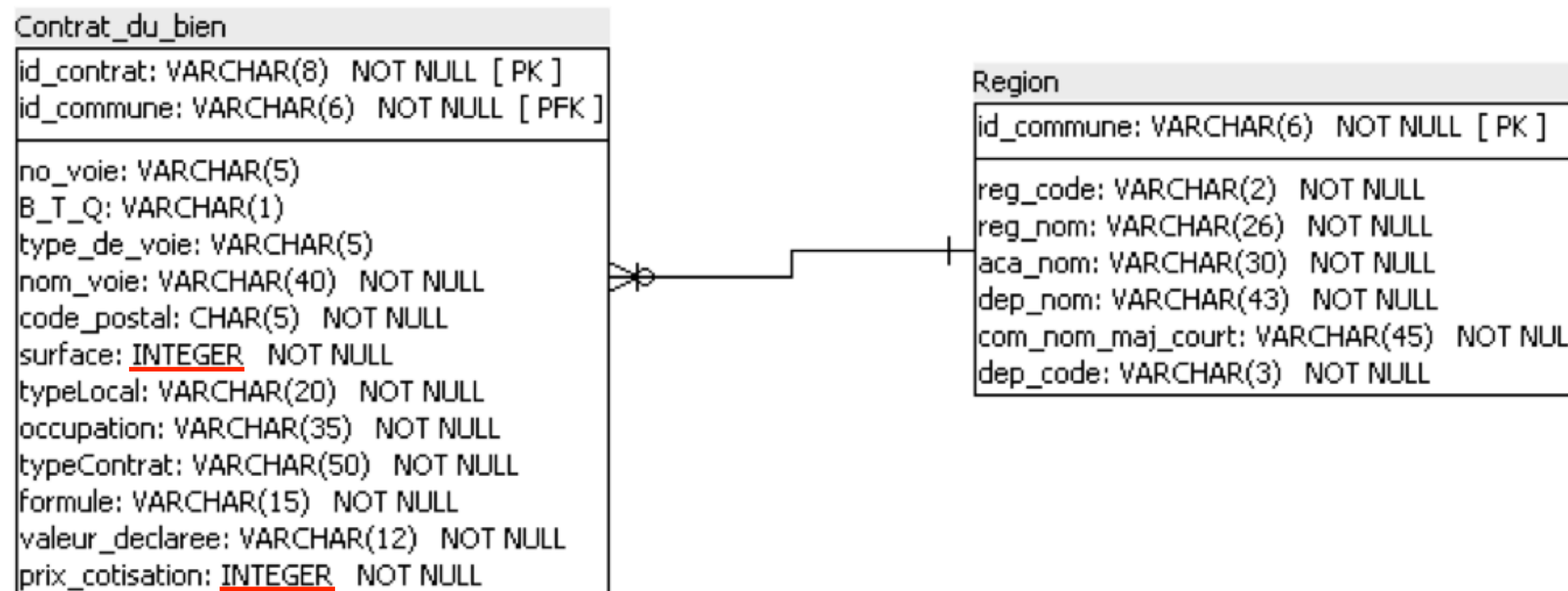
Une analyse initiale basée sur un modèle conceptuel UML a conduit à envisager une structure à trois tables pour mieux représenter certaines relations et détails des données.

Cependant, la structure finale a été ajustée à deux tables, conformément aux exigences du projet, afin de respecter les consignes et éviter des modifications importantes des fichiers de données fournies.

## 2.1 Modèle conceptuel UML à trois tables



## 2.2 Schéma relationnel normalisé 3NF à deux tables



# 3. Code SQL pour créer les tables

Le schéma relationnel diffère légèrement en raison des limites de *SQL Power Architect*, qui n'accepte pas *MEDIUMINT UNSIGNED* ni certaines contraintes, ajustées manuellement dans le code SQL final.

```
1 CREATE TABLE Region (  
2     id_commune VARCHAR(6) NOT NULL UNIQUE,  
3     reg_code CHAR(2) NOT NULL,  
4     reg_nom VARCHAR(26) NOT NULL,  
5     aca_nom VARCHAR(30) NOT NULL,  
6     dep_nom VARCHAR(43) NOT NULL,  
7     com_nom_maj_court VARCHAR(45) NOT NULL,  
8     dep_code VARCHAR(3) NOT NULL,  
9     PRIMARY KEY (id_commune)  
10 );  
11
```

```
12 CREATE TABLE Contrat_du_bien (  
13     id_contrat VARCHAR(8) NOT NULL UNIQUE,  
14     id_commune VARCHAR(6) NOT NULL,  
15     no_voie VARCHAR(5),  
16     B_T_Q VARCHAR(1),  
17     type_de_voie VARCHAR(5),  
18     nom_voie VARCHAR(40) NOT NULL,  
19     code_postal CHAR(5) NOT NULL,  
20     surface MEDIUMINT UNSIGNED NOT NULL,  
21     typeLocal VARCHAR(20) NOT NULL,  
22     occupation VARCHAR(35) NOT NULL,  
23     typeContrat VARCHAR(50) NOT NULL,  
24     formule VARCHAR(15) NOT NULL,  
25     valeur_declaree VARCHAR(12) NOT NULL,  
26     prix_cotisation MEDIUMINT UNSIGNED NOT NULL,  
27     PRIMARY KEY (id_contrat, id_commune)  
28 );  
29  
30 ALTER TABLE Contrat_du_bien ADD CONSTRAINT region_contrat_du_bien_fk  
31 FOREIGN KEY (id_commune)  
32 REFERENCES Region (id_commune)  
33 ON DELETE CASCADE  
34 ON UPDATE CASCADE;
```

# 4. Chargement de la base de données

## 4.1 Système choisi

*MySQL* et *MySQL WORKBENCH*

## 4.2 Correction des incohérences

Un problème a été détecté dans les données initiales : la clé primaire étrangère (*code\_dep\_code\_commune*) ne correspondait pas à la clé primaire de la table "REGION" dans trois cas spécifiques (97434, 97460, 97470).

À la place du code commune, le code postal avait été inséré par erreur. Cette incohérence empêchait l'insertion des données.

Pour résoudre ce problème, les valeurs ont été corrigées manuellement.

## 4.3 Chargement des données

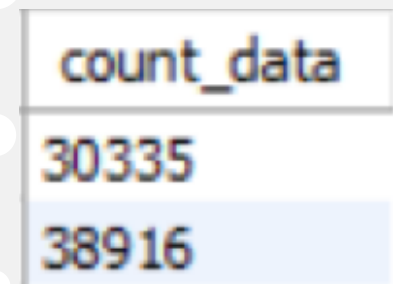
Les fichiers CSV ont été chargés dans les tables "contrat\_du\_bien" (30 335 lignes) et "region" (38 916 lignes),

```
Query OK, 30335 rows affected (1.62 sec)  
Records: 30335 Deleted: 0 Skipped: 0 Warnings: 0
```

```
Query OK, 38916 rows affected (1.28 sec)  
Records: 38916 Deleted: 0 Skipped: 0 Warnings: 0
```

puis leur insertion complète a été vérifiée, avec le code SQL suivant:

```
1 SELECT count(DISTINCT id_contrat)  
2 AS count_data FROM contrat_du_bien  
3 UNION  
4 SELECT count(DISTINCT id_commune)  
5 FROM region;
```



count_data
30335
38916

# 5. Requêtes SQL

Avec MySQL WORKBENCH

## 1 Comprendre la requête

*Quelle est la surface moyenne des contrats à Paris ?*

*Lister les numéros de contrats avec leur surface pour la commune de Caen.*

*Quels sont les 5 contrats qui ont les surfaces les plus élevées ?*

## 2 Traitement du code SQL

```
1 SELECT round(avg(c.surface)) AS Surface_Moyenne_m2
2 FROM contrat_du_bien c
3 JOIN region r ON c.id_commune = r.id_commune
4 WHERE r.com_nom_maj_court LIKE '%Paris %'
5 AND r.dep_nom = 'Paris';
```

```
1 SELECT c.id_contrat, c.surface AS surface_m2,
2 r.com_nom_maj_court
3 FROM contrat_du_bien c
4 JOIN region r ON c.id_commune = r.id_commune
5 WHERE lower(r.com_nom_maj_court) = 'caen';
```

```
1 SELECT id_contrat, surface AS surface_m2
2 FROM contrat_du_bien
3 ORDER BY surface DESC
4 LIMIT 5;
```

## 3 Visualisation du résultat sur MySQL

Surface_Moyenne_m2
52

id_contrat	surface_m2	com_nom_maj_court
103791	35	CAEN
103792	99	CAEN
103793	40	CAEN
103794	20	CAEN

id_contrat	surface_m2
104211	815
105463	742
130878	595
100822	570
109872	559

# 5. Requêtes SQL

## Autres Exemples

```
1 SELECT round(avg(prix_cotisation))
2 AS prix_moyen_par_mois
3 FROM contrat_du_bien;
```

Le prix moyen des cotisations mensuels est de 19 euros

```
1 SELECT valeur_declaree,
2 count(valeur_declaree)
3 FROM contrat_du_bien
4 GROUP BY valeur_declaree
5 ORDER BY no_contrats DESC;
```

La majorité des biens est déclarée pour une valeur inférieure à 25 000 euros

Seulement 104 biens ont déclaré une valeur supérieure à 100 000 euros

```
1 SELECT DISTINCT r.reg_code,
2 r.reg_nom AS Region,
3 count(c.id_contrat) AS no_contrats
4 FROM contrat_du_bien c
5 JOIN region r ON c.id_commune = r.id_commune
6 GROUP BY r.reg_code, r.reg_nom
7 ORDER BY no_contrats;
```

L'Île de France est la région avec le plus de contrats.

La région avec le moins de contrats est la Réunion

```
1 SELECT r.dep_code, r.dep_nom,
2 round(avg(c.prix_cotisation))
3 AS Prix_moyen
4 FROM contrat_du_bien c
5 JOIN region r
6 ON c.id_commune = r.id_commune
7 GROUP BY r.dep_code, r.dep_nom
8 ORDER BY Prix_moyen DESC LIMIT 10;
```

La cotisation moyenne la plus chère est à Paris: 36 euros

```
1 SELECT r.reg_nom, c.formule,
2 count(c.formule) AS no_contrats
3 FROM contrat_du_bien c
4 JOIN region r ON c.id_commune = r.id_commune
5 WHERE lower(r.reg_nom) LIKE '%pays%loire%'
6 AND lower(c.formule) = 'integrale'
7 GROUP BY r.reg_nom, formule;
```

Les contrats avec formule intégrale des Pays de la Loire sont 589

```
1 SELECT c.id_contrat, c.typeContrat,
2 c.formule, c.typeLocal, r.dep_code
3 FROM contrat_du_bien c
4 JOIN region r ON c.id_commune = r.id_commune
5 WHERE lower(c.typeLocal) LIKE '%maison%'
6 AND r.dep_code = '71'
7 ORDER BY c.id_contrat;
```

Dans le département 71, il y a 4 contrats pour les maisons, répartis équitablement entre les formules classique et intégrale.

# 6. Contraintes du projet

- — ○ **colonne dep\_nom\_num redondante**
- — ○ **Modèle conceptuel des données qui conduit sur un modèle à trois tables, à la place de deux**
- — ○ **3 données corrompues sur la clé étrangère empêchent l'insertion des données**
- — ○ **SQL Power Architect limite l'insertion de certains types de données**
- — ○ **Requête 2 égale à la requête 8**

# 7. Solutions proposées

- — ○ **colonne dep\_nom\_num pas prise en compte**

- — ○ **Retour au modèle original de deux tables du projet guidé, pour respecter les consignes**

- — ○ **3 données corrigées à la main en recherchant le bon code\_dep\_code\_commune**

- — ○ **Modifications souhaitées fait à la main directement dans le code SQL**

- — ○ **La requête 2 a été changée pour la différencier de la 8**