

Projet Python avancé

Analyse du marché immobilier
à partir de données issues du web scraping

Sommaire

Page 2 - Introduction

Page 3 - Collecte des données — **scraper.py**

Page 4 - Exploration des données — **exploration.ipynb**

Page 5 - Nettoyage et géocodage — **clean_data.py**

Page 6 - Synthèse de la méthodologie

Page 7 - Application web interactive — **dashboard.py**

Page 8 - Conclusion générale

INTRODUCTION

Contexte

Le marché immobilier français présente d'importantes variations de prix selon les territoires, les types de biens et les localisations. Comprendre ces disparités nécessite l'accès à des données fiables et leur analyse structurée.

Source des données

Nous avons collecté des annonces immobilières (maisons et appartements) à partir du site :

immobilier.notaires.fr

reconnu pour la qualité de ses informations.

Le scraping couvre 22 départements représentatifs du territoire français :

01, 06, 13, 21, 29, 30, 31, 33, 34, 35, 38, 42, 44, 49, 51, 59, 62, 69, 74, 75, 76, 83

Architecture du projet

Le projet s'articule autour de 4 composants Python :

<i>Fichier</i>	<i>Rôle</i>
scraper.py	Collecte des annonces
exploration.ipynb	Exploration et validation
clean_data.py	Nettoyage et géocodage
dashboard.py	Application streamlit interactive

Objectif

Proposer un outil simple, visuel et interactif permettant d'explorer les tendances du marché immobilier : carte des biens, analyses par territoire, comparaisons de prix au m².

Source : immobilier.notaires.fr

Celui-ci reposant en partie sur un chargement dynamique des pages, l'utilisation de Selenium a été retenue afin d'accéder au contenu HTML complet des annonces.

Une fois les pages chargées, le contenu est analysé à l'aide de BeautifulSoup afin d'extraire les informations essentielles de chaque annonce (titre, description, prix, surface, nombre de pièces, localisation et type de bien).

Le scraping est réalisé de manière automatisée sur plusieurs pages et départements, puis les données collectées sont exportées dans un fichier CSV pour les étapes suivantes du projet.

Extrait de scraper.py

```
12 # --- 1. Configuration du Projet ---
13 # Liste des départements à scraper
14 departements = ['75','01','13','69','31','06','44','34','33','59','83','42','76','21','38','35','51','49','30','29','74','62']
15 pages_to_scrape = 5
16 base_url = "https://www.immobilier.notaires.fr/fr/annonces-immobilieres-liste"
17
18 # Liste pour stocker tous les résultats
19 all_data = []
20
21 # --- 2. Configuration de Selenium (Mode Headless) ---
22 chrome_options = Options()
23 # Mode Headless (sans interface graphique) pour la rapidité et la stabilité
24 chrome_options.add_argument("--headless")
25 chrome_options.add_argument("--no-sandbox")
26 chrome_options.add_argument("--disable-dev-shm-usage")
27 # User-Agent pour éviter d'être bloqué
28 chrome_options.add_argument(
29     "user-agent=Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/91.0.4472.124 Safari/537.36")
30
```

Cet extrait présente le périmètre du scraping (départements et pagination) ainsi que la configuration de Selenium en mode headless pour charger les pages dynamiques du site et limiter les risques de blocage.

```
67 # --- Niveau 2 : Bouclier sur chaque lien pour aller chercher le DETAIL ---
68 for link in liens_annonces:
69     try:
70         driver.get(link)
71         # Pause pour charger la page de détail dynamique
72         time.sleep(random.uniform(2, 4))
73
74         soup_detail = BeautifulSoup(driver.page_source, 'html.parser')
```

Pour chaque annonce, le script accède à la page de détail via le navigateur automatisé. Un temps d'attente est appliqué pour assurer le chargement correct de la page, puis le contenu HTML est analysé avec BeautifulSoup afin d'accéder à la structure de l'annonce.

Les informations essentielles de chaque annonce sont extraites à partir de la structure HTML de la page, notamment le titre, le type de bien, la localisation et la description lorsque celle-ci est disponible.

```
135 # --- Sauvegarde des données ---
136 all_data.append({
137     'Departement': dept,
138     'URL': link,
139     'Titre': titre,
140     'Type_Bien': type_de_bien,
141     'Localisation': location,
142     'Prix': prix,
143     'Surface_m2': surface,
144     'Nb_Pieces': nb_pieces,
145     'Description': description
146 })
```

Les données extraites sont regroupées dans une structure intermédiaire afin de constituer un ensemble homogène prêt à être exporté et analysé.

```
157 df = pd.DataFrame(all_data)
158 df.to_csv('annonces_completes_notaires.csv', index=False, encoding='utf-8-sig')
```

Les données collectées sont structurées dans un DataFrame pandas puis exportées dans un fichier CSV, qui servira de base aux étapes de nettoyage, d'analyse et de visualisation.

Rôle du notebook

Il sert à comprendre la structure, évaluer la qualité des données et valider les règles de nettoyage.

Visualisation des premières annonces collectées (prix, surface, localisation, description).

1. Aperçu général des données

df.head(3)												
[10]	✓ 0.0s											
...	Département		URL		Titre		Type Bien	Localisation	Prix	Surface_m2	Nb_Pieces	Description
0	75	https://www.immobilier.notaires.fr/fr/annonce-...	Achat :Appartement T4- Paris 8 - Paris (75)		Appartement T4		- Paris 8 - Paris (75)	955000	72.12	4	Notaire - F4 Rue Berruyer PARIS Villème. Au coe...	
1	75	https://www.immobilier.notaires.fr/fr/annonce-...	Achat :Appartement T2- Paris 16 - Paris (75)		Appartement T2		- Paris 16 - Paris (75)	745500	70.00	2	Notaire vend un appartement T2 de 70.58m² Carr...	
2	75	https://www.immobilier.notaires.fr/fr/annonce-...	Achat :Appartement T4- Paris 15 - Paris (75)		Appartement T4		- Paris 15 - Paris (75)	600000	74.15	4	Exclusivité, au sein d'un bel immeuble d'angle...	

Extrait de : exploration.ipynb

Analyse des valeurs manquantes (nombre et pourcentage)

2.2 Valeurs manquantes

```

missing = df.isna().sum()
missing_pct = (missing / len(df) * 100).round(1)

pd.DataFrame({
    'Manquentes': missing,
    'Pourcentage': missing_pct
}))

```

✓ 0.0s

	Manquentes	Pourcentage
Departement	0	0.0
URL	0	0.0
Titre	0	0.0
Type_Bien	0	0.0
Localisation	0	0.0
Prix	0	0.0
Surface_m2	150	11.6
Nb_Pieces	0	0.0
Description	0	0.0

150 valeurs manquantes sur Surface m2 (11.6%)

→ À traiter dans `clean_data.py`

Analyse globale du jeu de données :

- volume d'annonces collectées,
- structure des colonnes et types de données,
- détection des doublons,
- identification des valeurs manquantes.

3. Traitement des surfaces manquantes

Hypothèse : La surface peut être extraite de la colonne `Description` via regex

3.1 Test du pattern regex principal

```
# Extraction de la surface depuis la liste à l'aide d'une expression régulière
pattern_m2 = r"^\(\s*\d+(\.\d+)\s*\)\s*(m²|m^2)$"

missing_surface = df["Surface_m2"].isna()

for i, row in missing_surface.head(3).iterrows():
    desc = row["description"][:200]
    match = re.search(pattern_m2, desc)
    print(f"Description: {desc}, {match}")

print("Surface extraite : {match.group(1)} if match else 'NON' {m²}")
```

Extraction des surfaces à l'aide d'expressions régulières (Regex) directement depuis la description des annonces, afin de corriger les valeurs manquantes.

Préparation de l'extraction de la ville pour le géocodage

Les informations de localisation étant fournies sous forme textuelle, cette étape vise à extraire le nom de la ville afin de préparer le géocodage (latitude / longitude) utilisé dans les analyses territoriales et la cartographie.

Cette phase d'exploration permet de formaliser précisément l'ensemble des règles de transformation nécessaires à la fiabilisation des données.

Ces règles seront ensuite implémentées dans le script `clean_data.py` afin de produire un jeu de données propre, structuré et prêt pour les analyses et la visualisation.

clean_data.py - Nettoyage et géocodage des données

Rôle du script clean_data.py

Le script clean_data.py transforme les données brutes issues du scraping en un jeu de données propre, structuré et exploitable, prêt à être utilisé dans l'application Streamlit.

Contrairement au notebook d'exploration, ce script constitue un pipeline de traitement automatisé, reproductible et orienté production.

Extrait de clean_data.py

```
1 """
2 clean_data.py - Script de nettoyage des données immobilières
3
4 Ce script transforme les données brutes du scraping en un fichier propre
5 prêt pour l'application Streamlit.
6
7 Entrée : data/annonces_raw.csv
8 Sortie : data/annonces_clean.csv
9
10 Transformations effectuées :
11 1. Extraction des surfaces manquantes depuis les descriptions
12 2. Suppression des lignes sans surface
13 3. Création de la colonne Ville
14 4. Création de la colonne prix_m2
15 5. Géocodage des villes (latitude, longitude)
16
17 """
```

Extraction intelligente des surfaces

```
87
88 # --- Priorité 1 : Surface Carrez (mention légale, très fiable) ---
89 carrez = re.search(r"(\d+(?:[.],)\d+)?\s*m[2]"s*(?:carrez|loi carrez)", desc, re.IGNORECASE)
90 if carrez:
91     return float(carrez.group(1).replace(',', '.'))
92
93 # --- Priorité 2 : Surface habitable ---
94 habitable = re.search(r"(\d+(?:[.],)\d+)?\s*m[2]"s*habitable?", desc, re.IGNORECASE)
95 if habitable:
96     return float(habitable.group(1).replace(',', '.'))
97
98 # --- Priorité 3 : Formulations courantes "de X m²" ---
99 # Capture : "appartement de 75 m²", "d'environ 80 m²", "d'une surface de 65 m²"
100 formulation = re.search(
101     r"(\d+(?:[.],)\d+)?\s*m[2]"s*(?:d'environ|d'une surface de)\s*m[2]",
102     desc, re.IGNORECASE
103 )
104 if formulation:
105     val = float(formulation.group(1).replace(',', '.'))
106     if val >= seuil_min:
107         return val
108
```

L'une des étapes clés du nettoyage consiste à corriger les surfaces manquantes à partir du texte des annonces.

Le script implémente une extraction hiérarchisée et robuste à l'aide d'expressions régulières.

La logique repose sur :

une priorisation des sources fiables (Surface Carrez, surface habitable), **des seuils adaptatifs selon le type de bien** (studios vs T2+), **l'évitement des pièges courants** (balcons, caves, terrasses).

Structuration des données géographiques

```
130 def extraire_ville(localisation: str) -> str:
131     """
132     Extraire le nom de la ville depuis la colonne Localisation.
133
134     Format attendu : "- Ville - Département (XX)"
135     Exemple : "- Paris 8 - Paris (75)" -> "Paris 8"
136
137     Args:
138         localisation: Chaîne de localisation brute
139
140     Returns:
141         Nom de la ville nettoyé
142     """
143     # Retirer le tiret initial et les espaces
144     texte = str(localisation).lstrip('- ').strip()
145
146     # Prendre la partie avant le second tiret
147     if '-' in texte:
148         return texte.split('- ')[0].strip()
149     return texte
150
```

Afin de permettre les analyses territoriales et la cartographie, le script extrait et normalise les informations de localisation.

Les transformations incluent :

extraction du nom de la ville depuis la colonne Localisation,

extraction du nom du département pour lever les ambiguïtés,

préparation des données pour le géocodage.

Géocodage des villes

```
234
235 def construire_query_geocodage(ville: str, nom_departement: str = "") -> str:
236
237     # Pattern pour détecter Paris/Lyon/Marseille + arrondissement
238     match = re.search(r"(paris|lyon|marseille)\s*(\d+)", str(ville).lower().strip())
239
240     if match:
241         nom_ville = match.group(1)
242         arrondissement = int(match.group(2))
243
244         # Codes postaux de base
245         base_cp = {
246             'paris': 75000,
247             'lyon': 69000,
248             'marseille': 13000
249         }
250
251         cp = base_cp.get(nom_ville, 0)
252         if cp > 0:
253             return f"{cp + arrondissement}, France"
254
255     # Cas standard : ajouter le département pour plus de précision
256     if nom_departement:
257         return f"{ville}, {nom_departement}, France"
258     else:
259         return f"{ville}, France"
260
```

Le script enrichit ensuite les données avec des coordonnées GPS :

géocodage via Nominatim (OpenStreetMap), **gestion spécifique des arrondissements** de Paris, Lyon et Marseille, **utilisation d'un cache** pour limiter les requêtes et respecter l'API.

Résultat final du pipeline

À l'issue du traitement :
les annonces sans surface exploitable sont supprimées,
le prix au m² est calculé,
les coordonnées géographiques sont ajoutées,

un fichier final est généré pour l'analyse et la visualisation : annonces_clean.csv

Synthèse de la méthodologie

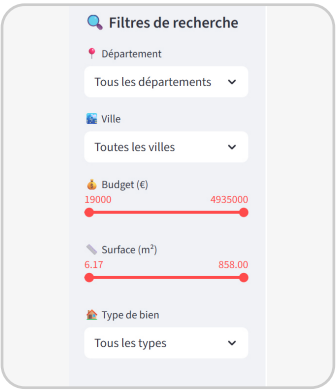
Le projet repose sur une chaîne de traitement complète et structurée des données immobilières. Il commence par le scraping automatisé des annonces depuis le site immobilier.notaires.fr, puis se poursuit par une phase d'exploration des données dans un notebook afin de mieux comprendre leur structure, leur qualité et les éventuels problèmes à corriger.

Les règles identifiées lors de cette exploration sont ensuite intégrées dans un script de nettoyage dédié, chargé de fiabiliser les données, de les enrichir (notamment sur le plan géographique) et de produire un jeu de données propre et cohérent. Cette organisation permet de bien séparer les rôles de chaque étape : exploration, nettoyage, enrichissement et préparation des données.

Les données finales ainsi obtenues peuvent alors être exploitées sereinement pour proposer des analyses statistiques et des visualisations interactives. La section suivante s'appuie sur ce jeu de données pour présenter une application Streamlit permettant d'explorer visuellement les tendances du marché immobilier.

dashboard.py - Application Streamlit

Les résultats du projet sont présentés à travers une application web interactive développée avec Streamlit, permettant d'explorer visuellement les données immobilières collectées et nettoyées. L'interface repose sur des filtres dynamiques et des visualisations interactives, offrant une analyse à la fois globale et détaillée du marché immobilier.



L'interface propose plusieurs sélecteurs interactifs permettant de filtrer les annonces par département, ville, type de bien (maison ou appartement), ainsi que par fourchette de prix et de surface à l'aide de curseurs.

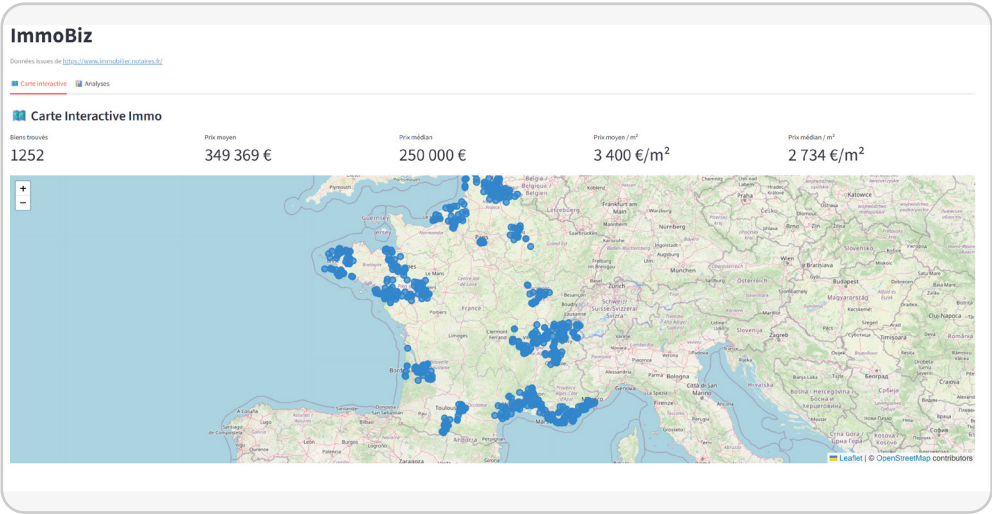
L'ensemble des visualisations se met à jour automatiquement en fonction des filtres sélectionnés, aussi bien pour la carte géolocalisée que pour les graphiques d'analyse.

L'application est organisée en deux onglets complémentaires.

Le premier onglet, **Carte interactive**, permet une exploration géographique des annonces, tandis que le second onglet, **Analyses**, propose des visualisations statistiques et comparatives des données.

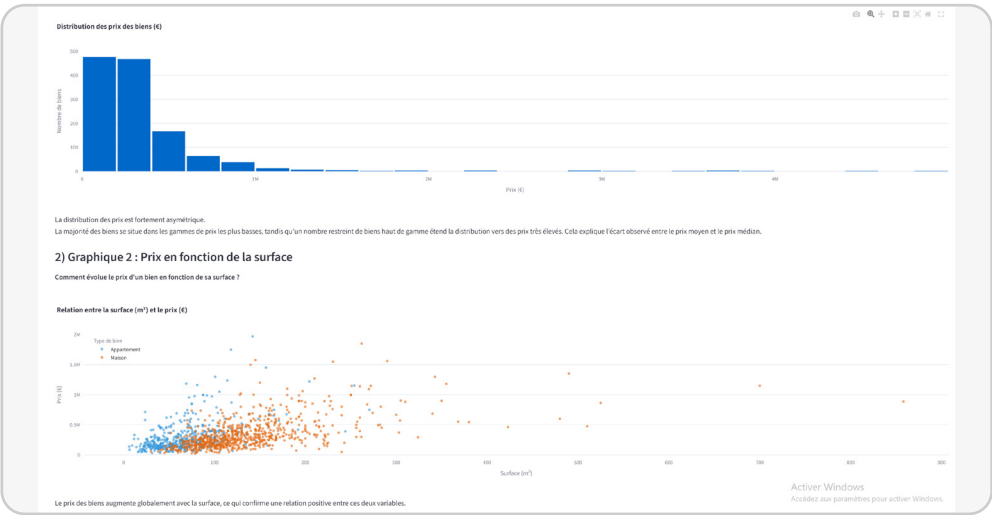
Carte interactive

Extrait de dashboard.py



Au-dessus de la carte, des indicateurs synthétiques permettent d'obtenir une lecture rapide du marché pour la sélection en cours. Sont notamment affichés le nombre de biens trouvés, le prix moyen, le prix médian, ainsi que le prix moyen et le prix médian au m².

Analyses



L'onglet Analyses regroupe plusieurs visualisations statistiques permettant d'explorer les données en profondeur. Il propose notamment la distribution des prix, la relation entre prix et surface, ainsi que des comparaisons du prix au m² par territoire, sous forme de graphiques et de distributions.

Conclusion

Ce projet a permis de mettre en place une chaîne complète de traitement et d'analyse de données immobilières, depuis la collecte automatisée des annonces jusqu'à leur visualisation au sein d'une application web interactive. Chaque étape a été pensée pour répondre à des problématiques concrètes de qualité des données, de structuration et de lisibilité des résultats.

La séparation claire entre le scraping, l'exploration, le nettoyage et la visualisation garantit une approche robuste, maintenable et évolutive. Les données finales obtenues sont fiables et suffisamment structurées pour permettre aussi bien des analyses globales que des explorations locales et détaillées du marché immobilier.

L'application Streamlit constitue l'aboutissement du projet en offrant une interface simple, interactive et accessible, permettant d'explorer les tendances du marché à travers des filtres dynamiques, une carte géolocalisée et des graphiques d'analyse. Ce projet illustre ainsi la capacité à concevoir une solution data complète, alliant traitement automatisé, analyse et restitution visuelle des données.