# Prédiction de l'environnement d'un quartier Stage de Master 2

Nelly Barret Université Claude Bernard Lyon 1

Encadrée par Fabien Duchateau et Franck Favetta

26 juin 2020









# Contexte du stage

#### Le projet Home in Love

- Objectif : aider à la recherche immobilière dans le cadre de la mobilité professionnelle
- Projet pluridisciplinaire débuté en 2017

#### Contexte scientifique

Comment qualifier simplement l'environnement d'un quartier ?

#### Objectif du stage

 Prédire l'environnement d'un quartier par apprentissage supervisé

## État de l'art

- Recommandation de logements [Yuan et al., 2013]
  - 3 critères : localisation, prix, unité urbaine
  - Utilisation d'une ontologie et du case-based reasoning
- Recommandation de quartiers [Liu et al., 2014]
  - Prise en compte du voisinage
  - Similarité entre quartiers (matrice) et régions (clustering)
- Comparaison manuelle de quartiers : datafrance.info
  - 5 critères : éducation, santé, services, commerces, loisirs
  - Classement des communes et visualisation cartographique

#### Positionnement

Recommandation de quartiers basée sur une expertise sociologique

Toward a user-oriented recommendation system for real estate websites. Exploiting geographical neighborhood characteristics for location recommendation.

## Quartier [Humain-Lamoure, 2007]

- Selon l'INSEE, un IRIS  $\rightarrow$  zone de 2000 à 5000 habitants
- En France, 50 000 IRIS

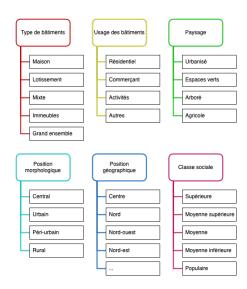


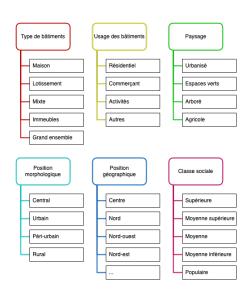
#### Indicateurs INSEE

Contexte

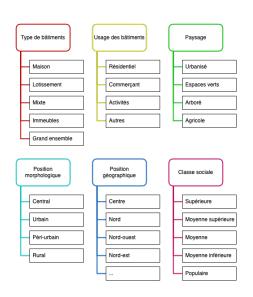
- 641 indicateurs
- Nombre de restaurants, population entre 18 et 25 ans, ...

- Résultat d'une analyse qualitative des sociologues sur 300 IRIS
- Description simplifiée d'un quartier par 6 critères





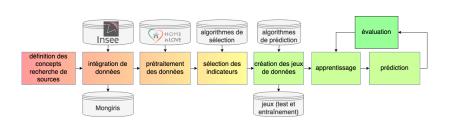






0000

# L'approche Predihood



#### Pré-traitement des données

#### IRIS expertisés

- Transformation des adresses en IRIS
- Traitement des valeurs inconnues : médiane des valeurs
- Traitement des valeurs erronées : semi-automatique

#### Indicateurs INSEE

- Traitement des valeurs inconnues : médiane des valeurs
- Normalisation : densité de population

# Représentativité des IRIS expertisés

Variable	Catégorie	Expertise	France
Position morphologique	rural	5%	68%
Daysaga	rural	17%	68%
Paysage	agricole	17/0	
Classe sociale	moyenne	82%	71%
	moyenne <sup>+</sup>	0270	
Type de bâtiments	collectif	68%	44%
Position géographique	nord, sud,	équitablement répartie	
Usage des bâtiments	nécessite une analyse particulière		

Tableau : L'analyse des IRIS expertisés montre que la plupart des variables d'environnement comportent un biais.

## Sélection des indicateurs

Objectif: générer des listes d'indicateurs utiles à la prédiction

## Étapes :

- Filtrage des indicateurs
  - 17 sont descriptifs (code postal, code IRIS...)
  - 208 sont trop spécifiques (nombre de courts de tennis couverts)
  - 59 sont non renseignés
- Sélection d'un sous-ensemble parmi 363 indicateurs
  - Combinaison d'algorithmes
  - Approche alternative : distribution des indicateurs

**Résultat** : plusieurs listes de k indicateurs pour chaque variable d'environnement v, notée  $L_{v}^{k}$ 

- Matrice de corrélation : supprimer les indicateurs 100% corrélés
- Algorithmes RF et ET : classer les indicateurs par importance
- Prise en compte de la diversité des catégories d'indicateurs

#### Algorithme 1 : Sélection des indicateurs pertinents pour la prédiction

```
Entrée : liste d'indicateurs \mathcal{I}, liste des variables d'environnement \mathcal{V}
     Sortie: listes d'indicateurs \hat{L}^k
 1 C \leftarrow \text{matriceCorrelation}(\mathcal{I}).\text{where}(\text{corr} = 1);
    \mathcal{I} \longleftarrow \mathcal{I} - C:
     for k \in [10, 20, 30, 40, 50, 75, 100] do
             for v \in \mathcal{V} do
                    L_{**} \longleftarrow \emptyset:
  5
                    F_{\cdot \cdot \cdot}^{ET} \leftarrow \text{top-k(ET.rank features}(\mathcal{I}), k):
  6
                    F_{ii}^{RF} \leftarrow \text{top-k(RF.rank features}(\mathcal{I}), k);
 7
                    F \longleftarrow F_{v}^{ET} \cup F_{v}^{RF};
                    for f \in \check{F} do
  9
                           p_f \leftarrow \operatorname{parent}(f);
10
                          if p_f \in F then
11
                                  p_f.score \leftarrow p_f.score + f.score;
12
13
14
```

# Interface cartographique





# Interface cartographique





# Interface de paramétrage





#### Parameters:

n\_estimators: 100 ; criterion: "gini" ; max\_depth: None ; min\_samples\_split: 2 ; min\_samples\_leaf: 1; min\_weight\_fraction\_leaf: 0 ; max\_features: "auto" ; max\_leaf\_nodes: None ; min\_impurity\_afercesses: 0 ; min\_impurity\_split: 1e-7 ; bootstrap: true ; oob\_score: false ; n\_lobs: None : con\_alons: 400 ; max\_samples: 4

Mean for this classifier: 52.14%

# Validation expérimentale

#### Niveau national

- Sélection de 5 algorithmes : Logistic Regression (LR), Random Forest (RF), K-Nearest Neighbours (KNN), Support Vector Classification (SVC) et AdaBoost (AB)
- Calcul de la précision pour chaque variable selon les 7 listes et les 5 algorithmes

#### Niveau communal

- Analyse qualitative des résultats prédits pour Lyon
- Effectuée par les sociologues

#### Résultats au niveau national

	LR	RF	KNN	SVC	AB
$\mathcal{I}$	52.9	64.5	59.3	<u>51.1</u>	55.6
L <sup>10</sup>	52.6	61.2	63.8	49.6	59.6
L <sup>20</sup>	55.9	64.1	63.0	49.6	56.6
L <sup>30</sup>	51.1	61.2	62.3	49.6	60.8
L <sup>40</sup>	<u>57.8</u>	63.0	60.8	49.2	56.3
L <sup>50</sup>	56.3	<u>64.9</u>	62.2	46.6	<u>61.1</u>
L <sup>75</sup>	50.7	63.4	60.8	51.1	58.2
L <sup>100</sup>	53.7	64.5	59.3	51.1	55.6

Tableau : Qualité de prédiction pour la variable usage

 Les listes améliorent la précision pour presque tous les algorithmes

#### Résultats au niveau national

Variable d'environnement	$\mathcal{I}$	$L^k$
Type de bâtiments	57%	60% (L <sup>20</sup> )
Usage des bâtiments	64%	65% ( <i>L</i> <sup>50</sup> )
Paysage	61%	63% (L <sup>20</sup> )
Classe sociale	51%	52%( <i>L</i> <sup>40</sup> )
Position géographique	34%	34% ( <i>I</i> )
Position morphologique	60%	61% ( <i>L</i> <sup>20,30,40</sup> )

Tableau : Qualité de prédiction pour l'algorithme Random Forest

- Random Forest est le meilleur algorithme
- La sélection permet une meilleure explicabilité des résultats

## Conclusion et perspectives

## Approche Predihood

- Algorithmes pour la prédiction de l'environnement (sélection des indicateurs et représentativité des données)
- Interfaces pour la visualisation des quartiers et le paramétrage d'algorithmes

#### **Perspectives**

- Sélectionner les indicateurs selon leur distribution
- Augmenter le jeu de données (avec les sociologues)
- Calculer la position géographique
- Intégrer de nouvelles sources de données (e.g. points d'intérêt)

Merci de votre attention !