

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 [YLKK13]**
 - 3 critères : localisation, prix, unité urbaine
 - Utilisation d'une ontologie et du *case-based reasoning*
- **Recommandation de quartiers [LWSM14]**
 - 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

Yuan, X. et al.: Toward a user-oriented recommendation system for real estate websites. *Information Systems* (2013).

Liu, Y. et al.: Exploiting geographical neighborhood characteristics for location recommendation. *Conference on Information and Knowledge Management* (2014).

État de l'art

- **Recommandation de logements** [YLKK13]
 - 3 critères : localisation, prix, unité urbaine
 - Utilisation d'une ontologie et du *case-based reasoning*
- **Recommandation de quartiers** [LWSM14]
 - 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

Recommandation de quartiers basée sur une expertise sociologique

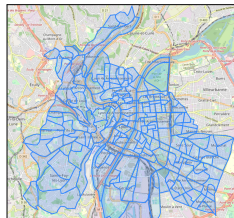
Yuan, X. et al.: Toward a user-oriented recommendation system for real estate websites. *Information Systems* (2013).

Liu, Y. et al.: Exploiting geographical neighborhood characteristics for location recommendation. *Conference on Information and Knowledge Management* (2014).

Prérequis

Quartier [HL07]

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



Indicateurs INSEE

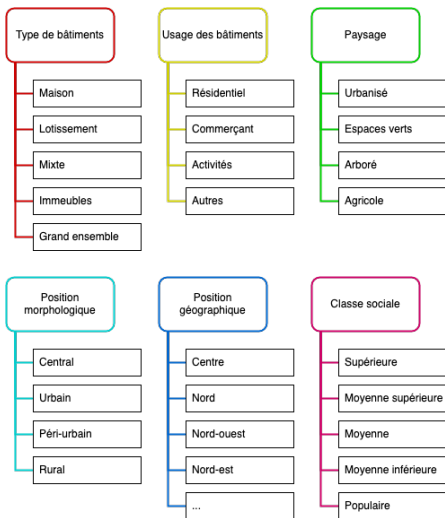
- 641 indicateurs
- Nombre de restaurants, population entre 18 et 25 ans, ...

Variables d'environnement

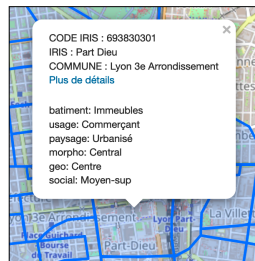
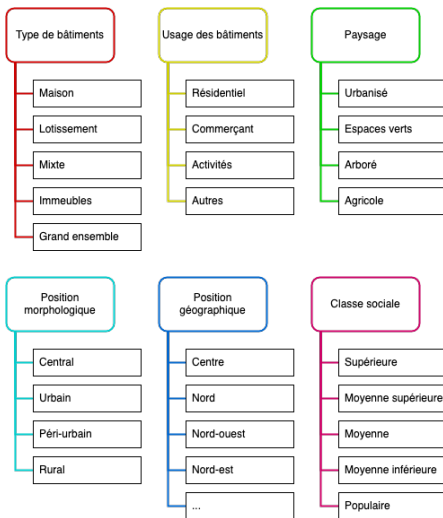
- Description simplifiée d'un quartier par 6 critères
- Résultat d'une analyse qualitative des sociologues sur 300 IRIS

Humain-Lamoure, A.L.: Le quartier comme objet en géographie. La Découverte (2007).

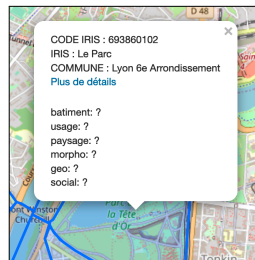
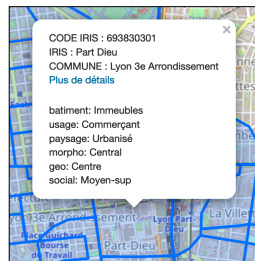
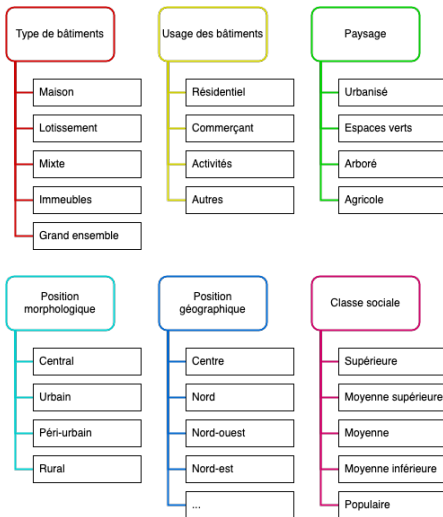
Variables d'environnement



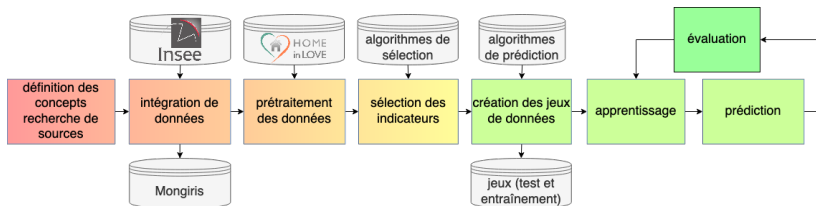
Variables d'environnement



Variables d'environnement



L'approche Predihood



Nguyen, G., et al.: Machine Learning and Deep Learning frameworks and libraries for large-scale data mining: a survey. AI Review (2019).

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%
Paysage	rural agricole	17%	68%
Classe sociale	moyenne moyenne ⁺	82%	71%
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 restants
 - Combinaison d'algorithmes
 - Approche alternative : distribution des indicateurs

Résultat : plusieurs listes de k indicateurs pour chaque variable d'environnement v , notée L_v^k

Sélection par combinaison d'algorithmes

- 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 L_v^k

```
1  $C \leftarrow \text{matriceCorrelation}(\mathcal{I}).\text{where}(\text{corr} = 1)$ ;  
2  $\mathcal{I} \leftarrow \mathcal{I} - C$ ;  
3 for  $k \in [10, 20, 30, 40, 50, 75, 100]$  do  
4   for  $v \in \mathcal{V}$  do  
5      $L_v \leftarrow \emptyset$ ;  
6      $F_v^{ET} \leftarrow \text{top-k}(\text{ET.rank\_features}(\mathcal{I}), k)$ ;  
7      $F_v^{RF} \leftarrow \text{top-k}(\text{RF.rank\_features}(\mathcal{I}), k)$ ;  
8      $F \leftarrow F_v^{ET} \cup F_v^{RF}$ ;  
9     for  $f \in F$  do  
10       $p_f \leftarrow \text{parent}(f)$ ;  
11      if  $p_f \in F$  then  
12         $p_f.\text{score} \leftarrow p_f.\text{score} + f.\text{score}$ ;  
13         $F \leftarrow F - \{f\}$ ;  
14    $L_v^k \leftarrow F$ ;
```

Interface Predihood

predihood
A tool for visualizing IRIS

203 iris found for query lyon.

Minimal zoom level to display IRIS automatically

12 (actual zoom level = 16)

Search by IRIS code
740560104

Search by IRIS name or city
lyon

Clear

LIRIS **LABEX IMU**

CCIS IRIS : 69080391
IRIS : Part Dieu
COMMUNE : Lyon 3e Arrondissement
Plus de détails
RandomForestClassifier
batiment: Immeubles (17)
usage: Commercial (897)
paysage: Urbain (17)
morpho: Centre (17)
geo: Centre (17)
social: Moyenne sup (97)

predihood

SELECT A CLASSIFIER
RandomForestClassifier

TUNE PARAMETERS
Classifier's parameters & Common parameters & TUNE DATASET

Train size: 80
Test size: 20
Remove outliers

Train, test and evaluate Abort

Results

	None	10	20	30	40	50	75	100	Mean
batiment	49.64%	50.76%	53.72%	53.73%	54.47%	53.37%	53.76%	54.11%	53.41%
usage	57.09%	58.57%	59.31%	57.46%	62.29%	58.57%	58.55%	60.06%	59.26%
paysage	54.45%	61.17%	57.49%	59.86%	57.44%	60.78%	56.69%	55.86%	58.40%
morpho	60.43%	61.18%	63.78%	61.91%	64.54%	66.65%	61.54%	64.51%	63.30%
geo	28.71%	29.45%	31.69%	34.31%	35.07%	33.57%	31.33%	31.68%	32.44%
social	42.19%	46.99%	48.89%	46.26%	46.64%	42.96%	44.76%	45.51%	46.00%

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_decrease: 0 ; min_impurity_split: 1e-7 ; bootstrap: true ; oob_score: false ; n_jobs: None ; random_state: None ; verbose: 0 ; warm_start: false ; class_weight: None ; ccp_alpha: 0.0 ; max_samples: None

Mean for this classifier: 52.14%

Validation expérimentale

Protocole

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^{10}	52.6	61.2	<u>63.8</u>	49.6	59.6
L^{20}	55.9	64.1	63.0	49.6	56.6
L^{30}	51.1	61.2	62.3	49.6	60.8
L^{40}	<u>57.8</u>	63.0	60.8	49.2	56.3
L^{50}	56.3	<u>64.9</u>	62.2	46.6	<u>61.1</u>
L^{75}	50.7	63.4	60.8	51.1	58.2
L^{100}	53.7	64.5	59.3	51.1	55.6

Tableau : Qualité de prédiction pour la variable *usage*

- Les listes améliorent la précision de la majorité des algorithmes

Résultats au niveau national

Variable d'environnement	\mathcal{I}	L^k
Type de bâtiments	57%	60% (L^{20})
Usage des bâtiments	64%	65% (L^{50})
Paysage	61%	63% (L^{20})
Classe sociale	51%	52% (L^{40})
Position géographique	34%	33% (L^{40})
Position morphologique	60%	61% ($L^{20,30,40}$)

Tableau : Qualité de prédiction pour l'algorithme Random Forest

- Random Forest obtient les meilleurs résultats
- 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)

Deuxième proposition

Distribution des indicateurs

Algorithme de sélection à partir de la distribution des indicateurs :

- Sélection des indicateurs discriminants
- Indépendant des variables d'environnement

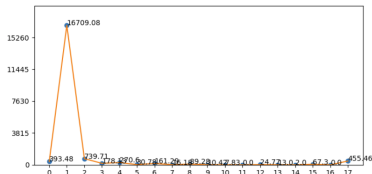


Figure: La Doua

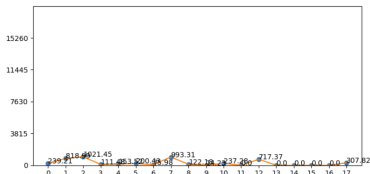


Figure: Saint-Cyr au Mont d'or