

Projet de Sciences des Données
Exploitation d'images satellites haute-résolution
pour la prevision d'indicateurs socio-économiques

YOUCEF - KACER

25 Aout 2016

Table des matières

Introduction	i
1 Images exploitées	1
2 Méthode d'apprentissage	3
3 Outils pour le stockage, le calcul et la visualisation	5

Introduction

Ce document présente le projet de sciences des données que je compte développer, dans le cadre de la validation du CES Data Scientist à Telecom-Paristech [TP14].

Ce projet consiste à exploiter des images satellitaires haute-résolution afin d'en extraire des indicateurs socio-économiques.

En effet, évaluer la densité démographique d'un pays, par exemple, peut représenter un coût non négligeable en terme de recensement. Or utiliser des images aériennes et leurs caractéristiques permettrait de prédire la population présente à moindre coût. En effet, les « edges » des routes et des bâtiments caractérisent les zones urbaines et donc les zones à forte population, alors que les champs et les forêts caractérisent des zones faiblement peuplées.

Chapitre 1

Images exploitées

Les images à exploiter proviennent du satellite Landsat 8 de la NASA et sont libres d'accès [Sur16]. Ce satellite scanne tout le globe terrestre tous les 16 jours. Ces images permettent donc non seulement d'étudier une zone à un moment donné mais aussi d'étudier son évolution sur une période donnée.

Ces images sont très riches dans la mesure où elles présentent en tout 11 canaux, 9 dans le visible et 2 dans l'infra-rouge, pour des résolutions allant de 15 à 60 mètres. Donc en plus des caractéristiques de formes, le niveau des images doit pouvoir nous renseigner sur la nature des matériaux et des objets présents au sein d'une zone (métal ou végétation par exemple), les canaux infra-rouges pourront très certainement quantifier la présence humaine.

Chapitre 2

Méthode d'apprentissage

L'idée serait de s'intéresser à une certaine zone (un pays par exemple), dont on aurait l'indicateur de densité de population (valeur à prédire) pour un grand ensemble de communes du pays.

On pourra alors récupérer plusieurs images satellitaires quadrillant ce pays, et attribuer à chacune d'elles sa valeur de densité de population (on doit pouvoir utiliser la latitude et la longitude d'une image pour retrouver la commune concernée).

Ainsi, on récupère un ensemble classique d'images labelisées par sa densité de population.

Ensuite, on pourra extraire des descripteurs de ces images (histogramme orienté du gradient [DT05], entre autre) auxquels on appliquera un algorithme de regression supervisé (la valeur à prédire, la densité de population, est plutôt continue que discrète).

On aurait donc un modèle de classification capable de prédire la densité de population d'une zone en fonction d'images satellites.

On pourra alors tester la généralisation du classifieur, en s'intéressant à d'autres pays.

Chapitre 3

Outils pour le stockage, le calcul et la visualisation

Les images seront stockées sur HDFS [Whi09] (soit en mode « single-node cluster », soit en mode « multi-node cluster » via le cluster de Telecom Paris-Tech), cela permettra d'extraire les descripteurs par Map/Reduce. Ces descripteurs seront aussi stockés de manière distribuée, via une table HBASE afin de pouvoir effectuer des requêtes et vérifier les valeurs calculées.

On utilisera alors la librairie de Machine Learning MLlib de Spark [MBY⁺16], dédiée à l'apprentissage sur données distribuées dans HDFS. Pour la régression, cette librairie ne permet cependant que les régressions linéaires (soit par moindres carrés linéaires, par Ridge ou par Lasso) et les régressions logistiques (pas de SVM régressif).

On testera toutes ces méthodes pour en analyser les erreurs en cross-validation. Enfin, pour la visualisation des résultats, on proposera une page html utilisant la librairie D3.js [Jai14] pour l'affichage de la carte de densité de population. Cette interface web sera aussi interactive que possible afin de permettre des zooms mais aussi l'observation de la densité de population sur une année antérieure.

Bibliographie

- [DT05] Navneet Dalal and Bill Triggs. Histograms of oriented gradients for human detection. pages 886–893, 2005.
- [Jai14] Abhijit Jain. Data visualization with the d3.js javascript library. *J. Comput. Sci. Coll.*, 30(2) :139–141, December 2014.
- [MBY⁺16] Xiangrui Meng, Joseph Bradley, Burak Yavuz, Evan Sparks, Shivararam Venkataraman, Davies Liu, Jeremy Freeman, DB Tsai, Manish Amde, Sean Owen, Doris Xin, Reynold Xin, Michael J. Franklin, Reza Zadeh, Matei Zaharia, and Ameet Talwalkar. Mllib : Machine learning in apache spark. *J. Mach. Learn. Res.*, 17(1) :1235–1241, January 2016.
- [Sur16] U.S. Geological Survey. <http://eros.usgs.gov/>, 2016.
- [TP14] Telecom-ParisTech. <http://www.telecom-evolution.fr/fr/formations-certifiantes/ces-data-scientist>, 2014.
- [Whi09] Tom White. *Hadoop : The Definitive Guide*. O’Reilly Media, Inc., 1st edition, 2009.