

Proposition de stage recherche en laboratoire 2019-2020

Titre : Transcription de documents anciens par apprentissage profond

Description du sujet :

Dans le cadre de leurs recherches, les historiens sont confrontés à la relecture de nombreux documents manuscrits anciens. Dès lors qu'il s'agit d'archives volumineuses (recensement, actes notariés, comptabilités, etc.), ce travail est fastidieux et souvent rendu complexe du fait de la présence d'écritures différentes.

Un premier projet exploratoire (2017/2018) a permis d'étudier la possibilité d'automatiser en partie ce travail de déchiffrement, en utilisant les acquis de la reconnaissance de caractères, après numérisation (photographie) des documents concernés.

Ce projet a permis de (i) photographier les documents correspondant à un corpus de taille relativement modeste (recensement de 1790 de la ville d'Arras), (ii) extraire et indexer manuellement quelques éléments pertinents d'informations manuscrites et (iii) étudier les méthodes d'apprentissage automatique susceptibles de fournir la meilleure réponse au problème.

Les premiers résultats de cette étude ont montré la faisabilité de la reconnaissance automatique de caractères anciens manuscrits, avec toutes les problématiques qui peuvent être associées à cette tâche : caractères disparus, abréviations, supports de mauvaise qualité, scribes multiples, etc.

En fonction des premiers résultats obtenus, il est maintenant envisagé une poursuite à ce travail, en visant à la fois des documents plus anciens, mais également plus volumineux, fournis par les historiens (cf exemple ci-après). Cette poursuite permettra d'étudier la mise au point d'une méthode d'analyse générique compatible avec des corpus d'origine et d'époques variées, ainsi que la définition d'outils efficaces pour l'historien, lors de la phase d'exploitation.

Ce module « end to end » sera basé sur une mise en place et une comparaison de structures d'apprentissage profond (Alexnet, ResNet, Autoencoder, codés en langage Python) pour obtenir la meilleure efficacité de transcription.

Objectifs :

Le projet consiste donc à appliquer la représentation d'attributs d'images de structures liées au « Deep Learning » à la transcription de documents anciens. Dans ce sens, les choix des structures deviennent prépondérants et ces structures seront évaluées suivant les caractéristiques classiques (précision, recall, ...). Pendant ce stage, une étude bibliographique poussée sera mise en œuvre (reconnaissance de caractères,

particularités des documents anciens des 17 et 18 éme siècles). Les programmes seront développés et intégrés dans le logiciel initial. L'étude sera étendue aux problèmes de sélection et extraction d'attributs, de classifications supervisées et semi-supervisées.

L'objectif final est de bâtir une méthode générique permettant d'appliquer ces méthodes à d'autres textes manuscrits de grand volume (Moyen-Age, ...), voir d'autres applications des SHS.

Références

Yahya Abbass, rapport de stage master TSI/STIP, ULCO-UL, 2017/18

Documents de la journée de restitution A2U du 17/10/2018.

1799
1770

Population de la
Ville d'Arras

Remise Canton M^{re} Dauchez C^{te}

Grande place 1799		Nombre
1	Les Carreaux de paupiers	11 Soier
	Enfants	1
	Domestiques	3
2	Le sieur Boyer et son Epouse	2
	deux filles	2
3	Le sieur Goussier	1
	Domestiques	3
	Joussier et sa femme	2
	1 Garçon	1
	2 filles	2
30	M ^{re} Billon	1 Locataire
	Fontaine de Guenard	1 id.
	Louise Boyer et son Epouse	2
29	un domestique	1
	Locataire	1
	Le sieur Mercier et son Epouse	2
28	Domestique	1
	Locataire	3
Casa.	Churet et son Epouse	2
	Enfants	3
	Garçon Boulanger	2
27	Le sieur Martel et son Ep.	2
	un enfant	1

Encadrant(s) : A. Bigand (MC, HdR)

email : bigand@univ-littoral.fr