
Katabase: À la recherche des manuscrits vendus

Simon Gabay^{*1,2}, Ljudmila Petkovic , Alexandre Bartz^{*3}, Matthias Gille Levenson^{*4},
and Lucie Rondeau Du Noyer

¹Université de Genève – Suisse

²Université de Neuchâtel – Suisse

³École nationale des chartes – École nationale des Chartes, Ecole Nationale des Chartes – France

⁴Histoire, Archéologie et Littératures des mondes chrétiens et musulmans médiévaux – École Normale Supérieure - Lyon – France

Résumé

Les marches de l'art, des livres ou des manuscrits sont tous relativement anciens, mais ne bénéficient cependant pas des mêmes outils pour la recherche. Des bases de données comme ArtPrice[1] existent pour les beaux-arts (peinture, sculpture...) et permettent de recenser les ventes. Des outils équivalents existent pour les livres anciens aux États-Unis [2], au Royaume-Uni [3], en Allemagne [4] ou en France [5].

Si certains index pour les ventes de livres anciens recensent bien les autographes, tous ne font pas [6], et les publications apparues tardivement ne reviennent pas sur les ventes passées. La documentation est donc disparate et fragmentaire, concernant une ressource de premier ordre pour les collectionneurs, mais aussi pour les philologues en quête de sources, les historiens du livre ou les adeptes de la Rezeptionsgeschichte qui peuvent s'intéresser aux prix ou aux noms des collectionneurs.

Enjeu

Si les principaux problèmes posés par la numérisation de catalogues comme la Revue des autographes sont connus, tout comme les enjeux de la détection d'un manuscrit revenant plusieurs fois sur le marché, parfois sous une forme fragmentaire [8], il nous a paru important d'améliorer notre algorithme de classification. Ce dernier doit en effet être implémentable dans une application web disponible en ligne, tout en étant capable de traiter de grandes quantités de données avec un maximum de précision.

L'enjeu est donc la conception d'un algorithme de classification assez précis pour reconnaître un même document, mais assez souple pour s'accommoder de variations plus ou moins importantes.

Stratégie

Afin d'accélérer le traitement de l'information et d'alléger le poids des fichiers mis en ligne, l'encodage XML-TEI, qui n'est qu'un format pivot, est abandonné au profit du JSON.

*Intervenant

Chaque fois que c'est possible une string est convertie en integer ou en float :

- Pour la longueur (`number_of_pages`) les documents incomplets sont ramenés à un nombre decimal ("une page et demie" → 1.5, "un quart de page" → 0.25. . .)

- Pour le format (`format`) le nombre de pliage est le chiffre retenu ("in-4o" → 4, "in-folio" → 1...)

- Pour la date (`date`) on utilise le format ISO YYYY-MM-DD ("3 mai 1645" → 1645-05-03, "septembre 1736" → 1736-09...)

- Le type de document (`term`) est converti en chiffre: ainsi L.a.s. (Lettre autographe signee) a le code 7, tandis que P.a.s. (Pièce autographe signee) a le code 2.

Les informations en JSON sont alors transformées pour faire une base de données orientée graphe, afin de faciliter la réconciliation des données.

Reconciliation

La transformation des données en graphe permet de simplifier le mécanisme de réconciliation : si chaque nœud représente un document vendu, il suffit d'ajouter une arête entre deux nœuds une fois atteint un certain degré de similarité.

Nous parlons de similarité et non d'identité stricte, car il n'est pas souhaitable de rechercher cette dernière : deux entrées différentes peuvent en effet renvoyer à un même document pour des raisons internes (deux fragments d'un même manuscrit) comme externes (une faute d'OCR). Il faut donc contourner ce problème via un algorithme de classification apte à gérer ces disparités.

A partir de la liste des documents vendus, chaque entrée est comparée avec les autres. Cette comparaison se fait sur la base des informations clés standardisées dans le fichier JSON : pour chacune de ces informations, un système de bonus/malus est appliqué. Si le score obtenu est supérieur à 0.6, alors les entrées sont considérées comme renvoyant à un même manuscrit.

La valeur de ces bonus/malus a été trouvée de manière expérimentale, sur la base de tests unitaires évaluant l'efficacité de l'algorithme. Ces valeurs sont susceptibles d'évoluer avec l'ajout de nouveaux manuscrits.

Applications

Une application en ligne [<https://katabase.herokuapp.com>] s'appuie sur les données en JSON pour l'affichage des catalogues, qui sont disponibles à la lecture, et sur l'algorithme de classification afin de proposer un double mode de lecture des résultats pour une requête dans la base : la liste des ventes et la liste des manuscrits vendus.

Les données disponibles proviennent pour l'instant presque essentiellement de catalogues de vente à prix marqués, publiés dans le dernier tiers du XIX^e siècle à Paris par Gabriel Charavay (le détail précis des catalogues numérisés est disponible dans l'application).

En faisant tourner l'algorithme sur ces données préliminaires, nous pouvons déjà offrir quelques premiers résultats. Nous avons pu définir un ratio de retour sur le marché des manuscrits : pour 44 333 manuscrits vendus, 3 364 ont été vendus au moins deux fois, soit environ 7,5%. À première vue, ces retours sur le marché sont marqués par une nette tendance baissière, notamment pour les manuscrits les plus chers, peu importe l'époque de l'auteur – la faible variation du franc à cette période et le de court laps de temps étudié permet par ailleurs une comparaison des prix malgré l'évolution du cours de la monnaie.

Recherches futures

Du point de vue philologique, la base de données ainsi que les capacités de classement développées pour l'application devraient permettre de retrouver plus aisément les sources des futures éditions, mais aussi de garantir l'authenticité des documents. Ces données devraient aussi être exploitables dans le cadre d'une approche distante du corpus afin d'étudier, par exemple, la construction du canon via la valeur marchande des auteurs.

Données et application

L'application web est disponible à l'adresse suivante : <https://katabase.herokuapp.com>.

Toutes les données utilisées pour ce projet sont disponibles en ligne à l'adresse suivante : <https://github.com/katabase>.

Bibliographie

<https://fr.artprice.com>.

American Book-Prices Current, New York, 1894/95-. ABPC tend avec le temps à repertorier de plus en plus de ventes européennes.

Book-Auction Records, London, 1902–1997 et *Book Prices Current*, London, 1887-1952.

Jahrbuch der Auktionspreise für Bücher, Handschriften und Autographen, Hamburg, 1950-. Au début *Jahrbuch der Auktionspreise für Bücher und Autographen*.

L'Argus mensuel du livre ancien et moderne, Promodis, Paris, 1981-,

Autograph Prices Current, London, 1914-1922

Simon Gabay, Lucie Rondeau Du Noyer et Mohamed Khemakhem, " Selling autograph manuscripts in 19th c. Paris : digitising the *Revue des Autographes* ", dans *Atti del IX Convegno Annuale AIUCD. La svolta inevitabile : sfide e prospettive per l'Informatica Umanistica*, Milan, Italy, 2020 (Quaderni di Umanistica Digitale), p. 113-118, url : <https://hal.archives-ouvertes.fr/hal-02388407>.

S. Gabay, L. Rondeau Du Noyer, Matthias Gille Levenson, Ljudmila Petkovic et Alexandre Bartz, " Quantifying the Unknown : How many manuscripts of the marquise de Sevigne still exist ? ", dans *Digital Humanities DH2020*, Ottawa, Canada, 2020 (DH2020 Book of Abstracts), url : <https://hal.archives-ouvertes.fr/hal-02898929> (visite le 23/11/2020).