



Analyse des logs de Gallica pour cartographier les profils des gallicanautes ^a

Philippe Chevallier (BNF), François Roueff (Télécom ParisTech)

a. D'après Nouvellet, Beaudouin, D'Alché-Buc, Prieur, and Roueff [2017]



Plan

Objectifs et données brutes

Mise en forme

Ce qu'il faut retenir



Objectifs et données brutes

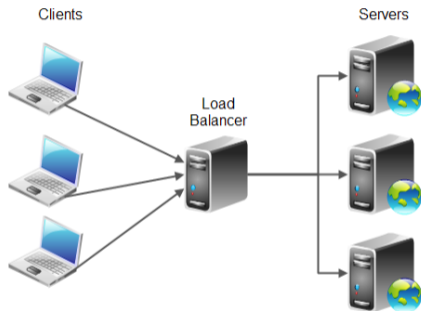
Mise en forme

Ce qu'il faut retenir

Objectifs

- **Objectif général** : améliorer le service de mise à disposition des ressources numériques aux *Gallicanauts*.
- **Contexte** : compléter des études et observations de terrain.
- **Contrainte** : exploration de toutes les données brutes disponibles (sur une durée donnée).
- **Approche** : modélisation des comportements d'utilisation des sites de la BnF.
- **Méthodologie** : classification non-supervisée des sessions vues comme séquences d'actions.
- **Outils** : Plateforme **Teralab**, programmation python, indexation par Elasticsearch.

Logs : architecture de Gallica/DataBnf



- Le Load Balancer redistribue les connexions entre les différents serveurs.
- Logs au niveau des 30 serveurs Apache.

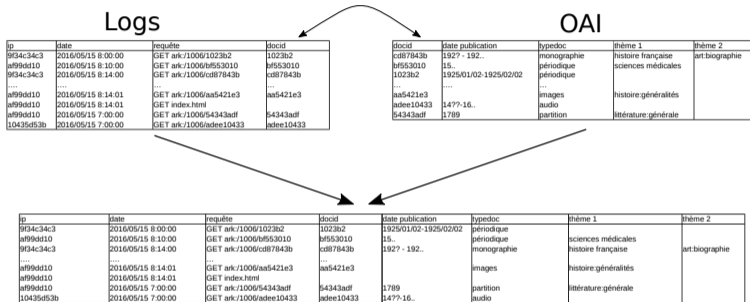
Logs : format

Format des logs :

- IP : **peut être partagée** par des utilisateurs d'une même institution
- pays, ville
- date
- requête : **ark/id unique**, page static, moteur de recherche, ...
- **referer** : indique le site de provenance
- user-agent : information sur le navigateur de l'utilisateur

```
##1becb38275ce3a5ee49b958c7e06d5b3##Algeria#null# - [12/Feb/2017:10:59:54 +0100] "GET /assets/static/javascripts/application/controllers/forms_controller.js HTTP/1.1" 200 2013 "http://gallica.bnf.fr/ark:/12148/bpt6k65089037.r=les%20ponts%20metalliques%20sur%20les%20
##1becb38275ce3a5ee49b958c7e06d5b3##Algeria#null# - [12/Feb/2017:10:59:54 +0100] "GET /assets/static/javascripts/application/controllers/contribution_controller.js HTTP/1.1" 200 1438 "http://gallica.bnf.fr/ark:/12148/bpt6k65089037.r=les%20ponts%20metalliques%20sur%20
##1becb38275ce3a5ee49b958c7e06d5b3##Algeria#null# - [12/Feb/2017:10:59:54 +0100] "GET /assets/static/javascripts/application/controllers/marquepage_controller.js HTTP/1.1" 200 906 "http://gallica.bnf.fr/ark:/12148/bpt6k65089037.r=les%20ponts%20metalliques%20sur%20les%20
##1becb38275ce3a5ee49b958c7e06d5b3##Algeria#null# - [12/Feb/2017:10:59:54 +0100] "GET /assets/static/javascripts/application/controllers/loader_controller.js HTTP/1.1" 200 1056 "http://gallica.bnf.fr/ark:/12148/bpt6k65089037.r=les%20ponts%20metalliques%20sur%20les%20
##1becb38275ce3a5ee49b958c7e06d5b3##Algeria#null# - [12/Feb/2017:10:59:54 +0100] "GET /assets/static/javascripts/application/controllers/issue-pagination_controller.js HTTP/1.1" 200 326 "http://gallica.bnf.fr/ark:/12148/bpt6k65089037.r=les%20ponts%20metalliques%20sur%20
##1becb38275ce3a5ee49b958c7e06d5b3##Algeria#null# - [12/Feb/2017:10:59:54 +0100] "GET /assets/static/javascripts/application/controllers/ariane-wire_controller.js HTTP/1.1" 200 318 "http://gallica.bnf.fr/ark:/12148/bpt6k65089037.r=les%20ponts%20metalliques%20sur%20les%20le
##1becb38275ce3a5ee49b958c7e06d5b3##Algeria#null# - [12/Feb/2017:10:59:54 +0100] "GET /assets/static/javascripts/application/controllers/viewer_controller.js HTTP/1.1" 200 1850 "http://gallica.bnf.fr/ark:/12148/bpt6k65089037.r=les%20ponts%20metalliques%20sur%20les%20
##1becb38275ce3a5ee49b958c7e06d5b3##Algeria#null# - [12/Feb/2017:10:59:54 +0100] "GET /assets/static/javascripts/main.js HTTP/1.1" 200 676 "http://gallica.bnf.fr/ark:/12148/bpt6k65089037.r=les%20ponts%20metalliques%20sur%20les%20chemins%20de%20fer%20algeriens?rk=6437
##1becb38275ce3a5ee49b958c7e06d5b3##Algeria#null# - [12/Feb/2017:10:59:54 +0100] "GET /assets/static/images/entete/appstore.png HTTP/1.1" 200 2519 "http://gallica.bnf.fr/ark:/12148/bpt6k65089037.r=les%20ponts%20metalliques%20sur%20les%20chemins%20de%20fer%20algeriens?rk
##1becb38275ce3a5ee49b958c7e06d5b3##Algeria#null# - [12/Feb/2017:10:59:54 +0100] "GET /assets/static/images/entete/gplay.png HTTP/1.1" 200 2589 "http://gallica.bnf.fr/ark:/12148/bpt6k65089037.r=les%20ponts%20metalliques%20sur%20les%20chemins%20de%20fer%20algeriens?rk
##1becb38275ce3a5ee49b958c7e06d5b3##Algeria#null# - [12/Feb/2017:10:59:54 +0100] "GET /assets/static/images/entete/facebook.png HTTP/1.1" 200 1352 "http://gallica.bnf.fr/ark:/12148/bpt6k65089037.r=les%20ponts%20metalliques%20sur%20les%20chemins%20de%20fer%20algeriens?rk
##1becb38275ce3a5ee49b958c7e06d5b3##Algeria#null# - [12/Feb/2017:10:59:54 +0100] "GET /assets/static/images/entete/twitter.png HTTP/1.1" 200 1524 "http://gallica.bnf.fr/ark:/12148/bpt6k65089037.r=les%20ponts%20metalliques%20sur%20les%20chemins%20de%20fer%20algeriens?rk
##1becb38275ce3a5ee49b958c7e06d5b3##Algeria#null# - [12/Feb/2017:10:59:54 +0100] "GET /assets/static/images/entete/pinterest.png HTTP/1.1" 200 1599 "http://gallica.bnf.fr/ark:/12148/bpt6k65089037.r=les%20ponts%20metalliques%20sur%20les%20chemins%20de%20fer%20algerien
##1becb38275ce3a5ee49b958c7e06d5b3##Algeria#null# - [12/Feb/2017:10:59:54 +0100] "GET /nbImage/perso/logo_header_1.png HTTP/1.1" 200 1159 "http://gallica.bnf.fr/ark:/12148/bpt6k65089037.r=les%20ponts%20metalliques%20sur%20les%20chemins%20de%20fer%20algeriens?rk=64378
##0226ae7c8fa9549382333c6d55a95##Thailand#null# - [12/Feb/2017:10:59:54 +0100] "GET /services/ajax/extract/ark:/12148/otv1b9034116k.r=Bugatti?rk=2639498;0 HTTP/1.0" 200 458 "http://gallica.bnf.fr/services/engine/search/sru/operation=search?retLevel=version=1.26star
##1becb38275ce3a5ee49b958c7e06d5b3##Algeria#null# - [12/Feb/2017:10:59:54 +0100] "GET /nbImage/perso/logo_header_2.png HTTP/1.1" 200 1372 "http://gallica.bnf.fr/ark:/12148/bpt6k65089037.r=les%20ponts%20metalliques%20sur%20les%20chemins%20de%20fer%20algeriens?rk=64378"
```

Enrichissement des données : OAI



Moissonnage et nettoyage de l'entrepôt OAI :

- titre, auteur.
- date de publication.
- catégories

Tri des données de log

- Une **session** est une séquence de requêtes faite par un unique usager durant la visite d'un site particulier.
- La sessionisation est l'étape permettant de trier l'ensemble des requêtes en sessions.
 - Séquence de requêtes présentant **la même IP**
 - Nouvelle session lorsque l'IP devient inactive pendant 60min¹
- Un robot internet est un logiciel qui réalise des tâches automatisées sur internet
 - "good bots" (robot d'indexation, droit d'auteur)
 - "bad bots" (robot moisonneur, botnet)

1. cf. enquête quali.



Objectifs et données brutes

Mise en forme

Ce qu'il faut retenir

Mise en forme/mise en “modèle”

- Mise en forme basique des données : collection de sessions, chacune vue comme une succession d'événements (plus ou moins riches) estampillés temporellement.
- Nécessité d'un **modèle** pour une classification non-supervisée permettant une **interprétation**²
 - Exemple 1 : suite ordonnée d'actions de l'utilisateur.
 - Exemple 1 enrichi : ajout d'une action basée sur le **referer**.
 - Exemple 2 : suite d'actions associées à des **durées**.

2. cf. enquête quali.

Exemple 1 : actions

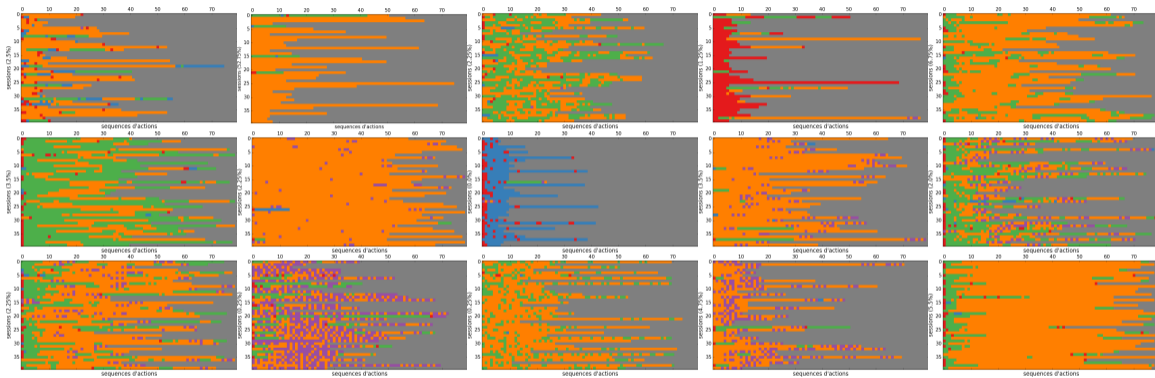
1. Navigation sur la page d'accueil : ■
2. Navigation sur les pages de médiations (collection + blog) : ■
3. Recherche Gallica : ■
4. Téléchargement d'un document : ■
5. Consultation locale dans l'interface Gallica (zoom, page suivante...) : ■

Une session de 10 actions

Représentation sous forme d'une "ligne" :



Exemple 1 : classification à 15 classes





Objectifs et données brutes

Mise en forme

Ce qu'il faut retenir

Importance de l'analyse des logs pour la connaissance des usages

- en **complément** (et non substitution) des autres méthodes
- prend en compte **tous les usages**, en particulier les plus difficiles à capter dans les enquêtes (usages "furtifs")
- permet une approche à la fois **fine et quantifiée** des parcours que le déclaratif rend imparfaitement
- ouvre des pistes d'**amélioration** concrètes en comparant le parcours idéal (pensé par le producteur) et le parcours réel (de l'utilisateur)

Quelle place pour l'analyse de traces d'usage à la BnF ?

- **limite des outils propriétaires** (mesure d'audience, e.réputation, etc.) : pas d'accès aux données brutes ; requêtes prédéfinies
- choix fort d'en faire un **objet de recherche** avec Télécom ParisTech : transparence des traitements / modèles évolutifs / interprétation collective
- nécessite un **dialogue** entre :
 - les producteurs des données (DSI de la BnF)
 - les experts de la collection numérique (bibliothécaires de la BnF)
 - les chercheurs en **science des données** (département IDS de Télécom)
 - les chercheurs en **sciences sociales** (département SES de Télécom)
- souhait de poursuivre ce type de traitements, en particulier pour mesurer l'impact des évolutions de l'interface sur les comportements



Bibliographie I

Adrien Nouvellet, Valérie Beaudouin, Florence D'Alché-Buc, Christophe Prieur, and François Roueff. Analyse des traces d'usage de gallica. Technical report, Télécom ParisTech, 2017.