

Historical big-data: modelization of strategies to analyze collections of documents

Camille Guerry

Supervisor : Bertrand Couâsnon, Aurélie Lemaitre, Sébastien Adam

SIFED : 07/06/2019

Introduction

Context of the project

ANR project: HBDEX - Exploitation of Big Data for the Historical Social Sciences application to Financial Data

Global aim: Understanding the long term dynamics of financial markets.

Collaboration with :

- DFIH team of the Paris School of Economics
- LITIS laboratory (Rouen)

Documents processing :

- Strategy for the recognition of collections
- Structure recognition
- Text recognition → LITIS

Introduction

Documents
characteristics

Collection
recognition
strategy

First
implementation

Conclusion

Global aim of the project

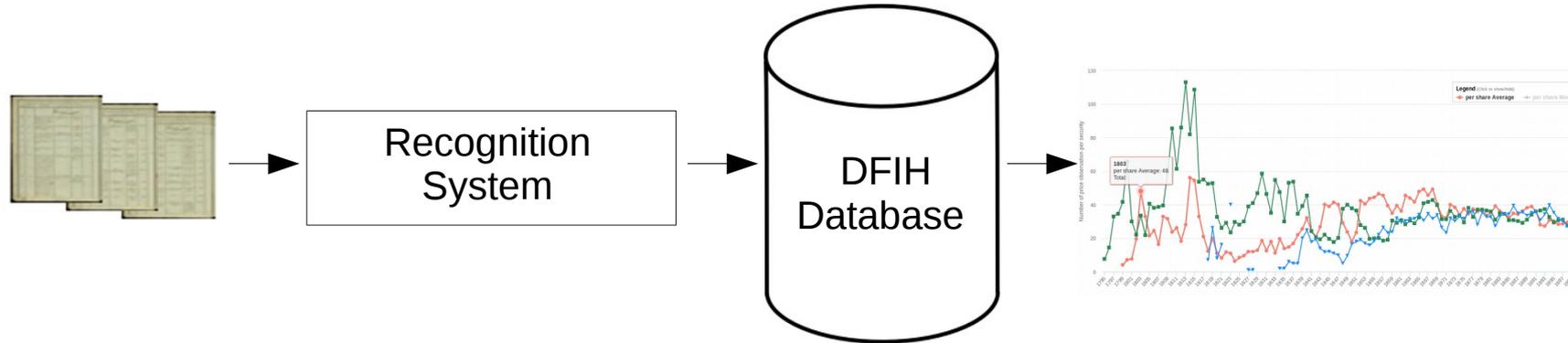
Introduction

Documents
characteristics

Collection
recognition
strategy

First
implementation

Conclusion



- Extraction of all information contained in the documents
- Limitation of human interaction
- Production of reliable data

The documents

- Kind of documents: **Stock exchange daily list**

- In the late 19th and early 20th centuries

- Printed Documents

- Different Markets

Stock exchange daily list:

- List of all stocks

- Tabular structure

Introduction

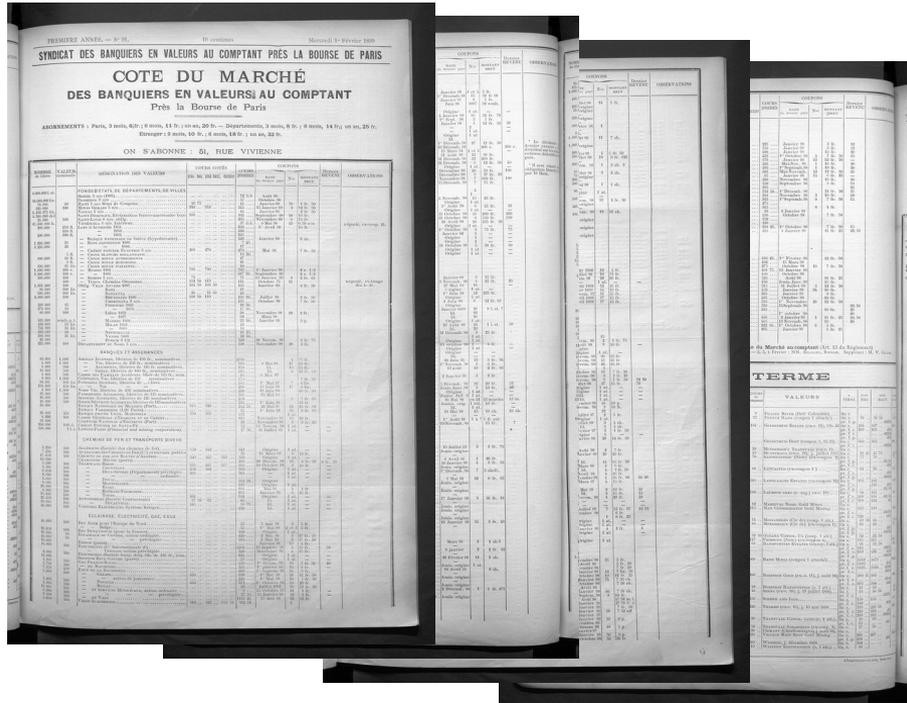
Documents characteristics

Collection recognition strategy

First implementation

Conclusion

La Coudisse (Paris unoficial market)



- Sequential data
- 1 day = several pages
- More than 140 000 images
 - 40 years (from 1899 to 1939)
 - between 4 and 15 pages/day

Work done on the Coullisse collection will be generic enough to be adapted to other collections

Introduction

Documents characteristics

Collection recognition strategy

First implementation

Conclusion

Presentation Content

Introduction

Document characteristics

Collection recognition strategy

First implementation

Conclusion

Introduction

Documents
characteristics

Collection
recognition
strategy

First
implementation

Conclusion

Introduction

**Documents
characteristics**

Collection
recognition
strategy

First
implementation

Conclusion

Documents characteristics

Extraction example

NOMBRE de titres	VALEUR nominale	DÉSIGNATION DES VALEURS	COURS COTÉS plus bas, plus haut, dernier	COURS précédents	COUPONS		
					DATE du dernier payé	Nos	MONTANT BRUT
		MÉTALLURGIE, FABRIQUES DE MACHINES					
5.000	500	ARIÈGE (Société métallurgique) ..	93	300	Janvier 98	3 et 4	5 fr.
35.000	100	ATELIERS FRANCO-RUSSES ..		92 50	1 ^{er} Décemb. 98	15	32 fr 60
5.000	500	— DU NORD DE LA FRANCE ..		595	Janvier 99	8	5 fr.
20.000	400	BOUHEY (Usines) ..		150	Janv 98	1897	32 rous.

Stock name : Ateliers du Nord de la France

Sector : MÉTALLURGIE, FABRIQUES DE MACHINES

Stock Quantity : Value : 5.000 Currency :

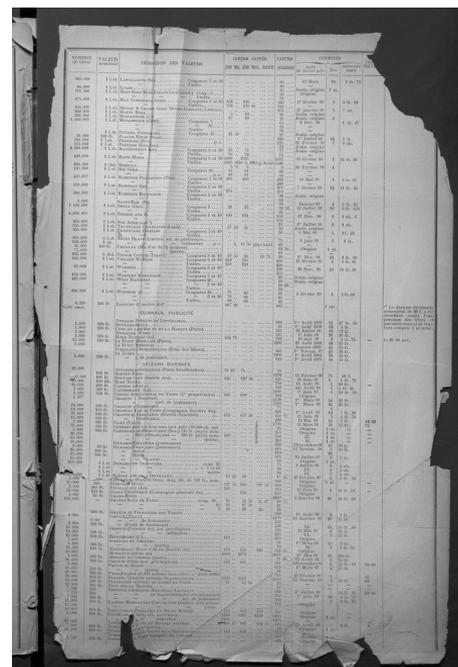
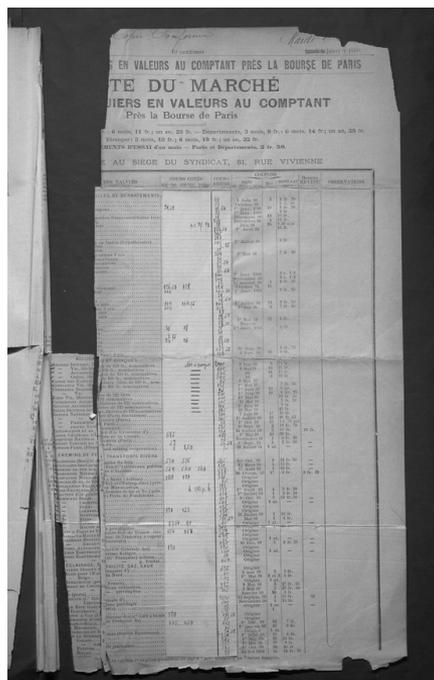
Stock Nominal : Value : 500 Currency : Comments :

Prices : Min : Max : Close : Previous: 595 Comments :

Dividend : Last payment date : 1898 - 12 - 01 Last_payment_date_precision : DAY
 Coupon_number : 15 Coupon_is_attached : N Gross value : 17.50
 Coupon_number_comments :

Difficulties

- Some documents are strongly deteriorated



Introduction

Documents characteristics

Collection recognition strategy

First implementation

Conclusion

Difficulties

- Some documents are strongly deteriorated
- Complex logical structure

65 ..	44 80	MEXICAIN int. 5 o/o j. avril 1899	fin e.....	44 90.	44 95	44 95
		(ex-coupon, 1 25 net)	pr. f. c. d/0f50	
			pr. f. c. d/0f25	
			pr. f. p. d/0f50	
			pr. f. p. d/0f25	

Introduction

Documents
characteristics

Collection
recognition
strategy

First
implementation

Conclusion

Difficulties

- Some documents are strongly deteriorated
- Complex logical structure
- Stock identification

Different cases:

- different name for the same stock
- close name for different stock
- specificity of an corporate action

Ex : 05/01/1899

05/01/1899

CONSTR. MÉCAN. DU MIDI DE LA RUSSIE, actions anc...
— — — — — nouv..
— — — — — 1/10 p. de Fond.

07/01/1899

CONSTR. MÉCAN. DU MIDI DE LA RUSSIE, actions
— — — — — 1/10 p. de Fond.

→ same id

Difficulties

- Some documents are strongly deteriorated
- Complex logical structure
- Stock identification

 **Use of the collection characteristics**

Introduction

**Documents
characteristics**

Collection
recognition
strategy

First
implementation

Conclusion

Collection characteristics

Structural stability of column

10/01/1899

NOMBRE de titres	VALEUR nominale	DÉSIGNATION DES VALEURS	COURS COTÉS plus bas, plus haut, dernier		CLOTURE précédente	COUPONS	
						DATE du dernier coupon	Montant et Spécification

16/01/1899

NOMBRE de titres	VALEUR nominale	DÉSIGNATION DES VALEURS	COURS COTÉS plus bas, plus haut, dernier		Derniers COURS	COUPONS			Dernier REVENU	OBSERVATIONS
						DATE du dernier payé	N ^{os}	MONTANT BRUT		

01/01/1905

NOMBRE DE TITRES	VALEUR NOMINALE	DÉSIGNATION DES VALEURS	COURS COTÉS plus bas plus haut		COURS précédents	COUPONS			DERNIER REVENU	OBSERVATIONS
						DATE du dernier payé	N ^{os}	MONTANT NET		

02/01/1914

MONTANT ou NOMBRE DE TITRES admissibles à la Cote	VALEUR NOMINALE	DÉSIGNATION DES VALEURS	COURS DU JOUR		COURS PRÉCÉDENTS de la semaine Plus bas Plus haut	RELEVÉ HEBDOMADAIRE des cours extrêmes du 22 au 27 Décembre 1913 ou antérieurement			COUPONS		
			PLUS BAS	PLUS HAUT		Plus bas	Plus haut	date	DATE du dernier payé	N ^{os}	MONTANT NET

28/12/1939

Vendredi 31 Mars 1939

Page 2

MONTANT ou NOMBRE DE TITRES		VALEUR nominale	COURS de compens ^{on}	COURS de REPORTS	DÉSIGNATION DES VALEURS	ÉCHÉANCES	MONTANT des PRIMES	COURS DU JOUR				CLOTURE précédente	COUPONS	
EMIS	ADMIS							Premier cours	Plus bas	Plus haut	Dernier cours		DATE DU DERNIER PAYÉ	MONTANT NET

Introduction

Documents characteristics

Collection recognition strategy

First implementation

Conclusion

Collection characteristics

Structural stability of sections

Introduction

Documents characteristics

Collection recognition strategy

First implementation

Conclusion

Sections layout

Sections names evolution :

Ordre d'apparition	1899 (07/01/1899)	1905 (21/07/1905)
1	FONDS D'ÉTATS, DE VILLES, DE DÉPARTEMENTS	FONDS D'ÉTATS, VILLES, DÉPARTEMENTS
2	ASSURANCES ET BANQUES	ASSURANCES & BANQUES
3	CHEMINS DE FER ET TRANSPORTS DIVERS	CHEMINS DE FER & TRANSPORTS
4	ÉCLAIRAGE, ÉLECTRICITÉ, GAZ, EAUX	ÉCLAIRAGE, ÉLECTRICITÉ, GAZ, EAUX
5	METALLURGIE, FABRIQUES DE MACHINES	INDUSTRIES TEXTILES, SOIES ARTIFICIELLES, TEINTURERIES
6	ACTIONS.— MINES DE CHARBON, CUIVRE, FER, ZINC, PLOMB, ARGENT, DIAMANTS.	MÉTALLURGIE, CONSTRONS MÉCANIQUES
7	MINES D'OR ET SOCIÉTÉS D'EXPLORATIONS	MINES DE CHARBONS, CUIVRE, FER, ZINC, PLOMB, ARGENT DIAMANTS
...

Collection characteristics

Structural stability of stock

16/02/1899

RAFFINERIE SAY
SUD-RUSSE, Soude.....
TAVERNES POUSSET ET ROYALE RÉUNIES.....
TURBIE-SUR-MER, Compagnie du Littoral.....

...

21/02/1899

RAFFINERIE SAY
SUD-RUSSE, Soude.....
TAVERNES POUSSET ET ROYALE RÉUNIES.....
TURBIE-SUR-MER, Compagnie du Littoral.....

23/02/1899

RAFFINERIE SAY
SUD-RUSSE, Soude.....
TAVERNES POUSSET ET ROYALE RÉUNIES.....
TERRAINS D'ARLES.....
TURBIE-SUR-MER, Compagnie du Littoral.....

“Terrain d’Arles”
is added

28/02/1899

RAFFINERIE SAY
SUD-RUSSE, Soude.....
TAVERNES POUSSET ET ROYALE RÉUNIES.....
TERRAINS D'ARLES.....
TURBIE-SUR-MER, Compagnie du Littoral.....

Introduction

Documents
characteristics

Collection
recognition
strategy

First
implementation

Conclusion

Difficulties and Challenges

Difficulties

- Some documents are strongly deteriorated
- Complex logical structure
- Stock identification

Challenges

- Use the structural stability of the collection and take the evolution into account
- Use economical rules that apply to the collection
- Propose a strategy usable for other collections

Introduction

**Documents
characteristics**

Collection
recognition
strategy

First
implementation

Conclusion

Collection recognition strategy

Introduction

Documents
characteristics

**Collection
recognition
strategy**

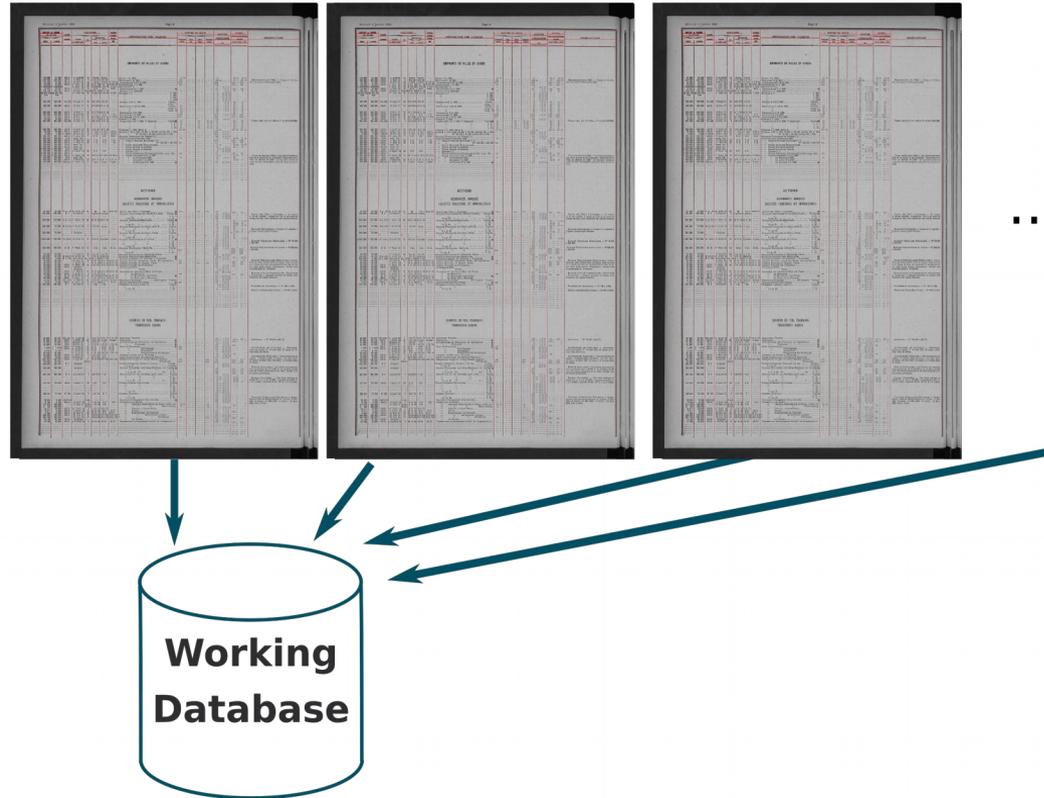
First
implementation

Conclusion

Iterative process

Structural Analysis

Columns



Introduction

Documents characteristics

Collection recognition strategy

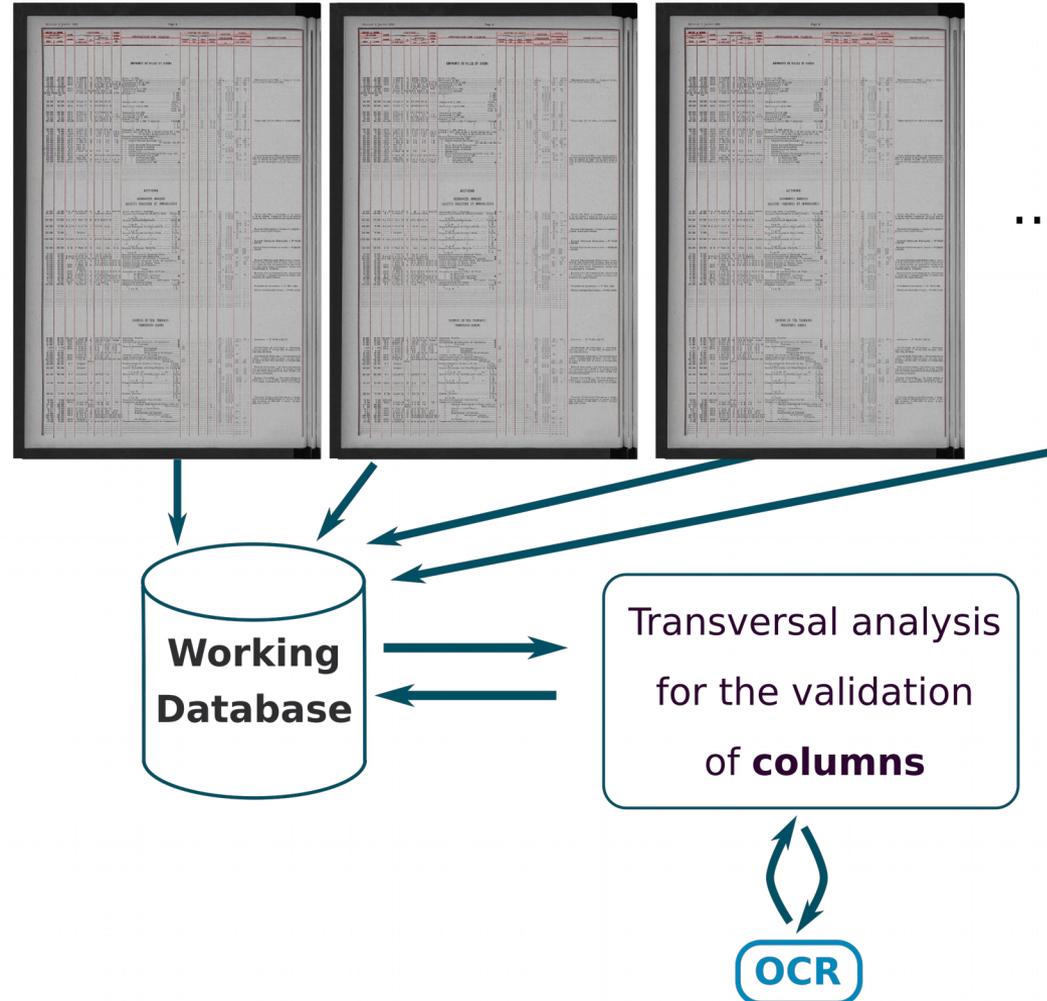
First implementation

Conclusion

Iterative process

Structural Analysis

Columns



Introduction

Documents characteristics

Collection recognition strategy

First implementation

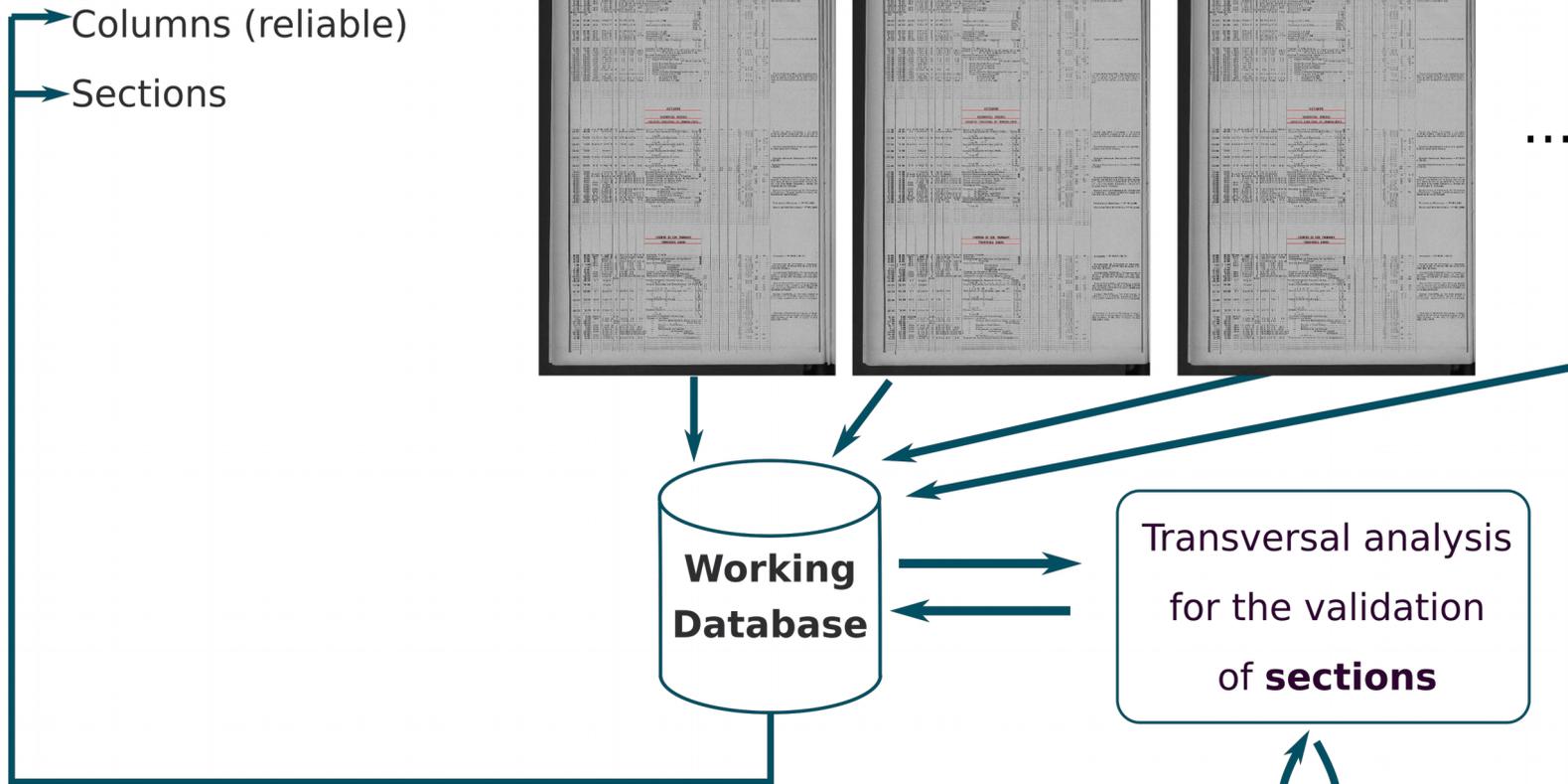
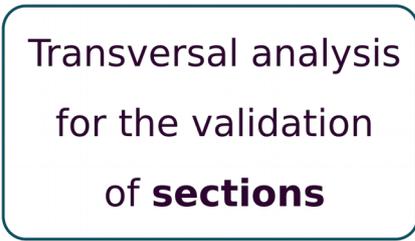
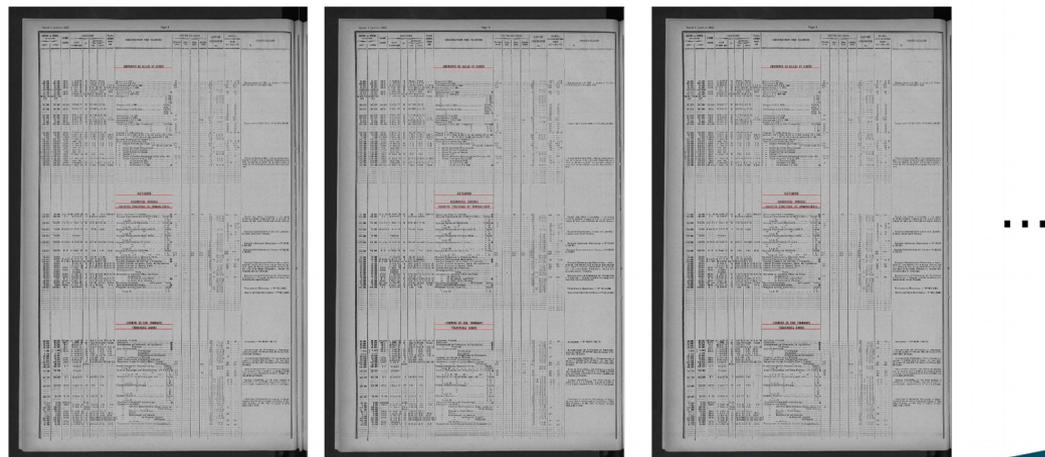
Conclusion

Iterative process

Structural Analysis

Columns (reliable)

Sections



Introduction

Documents characteristics

Collection recognition strategy

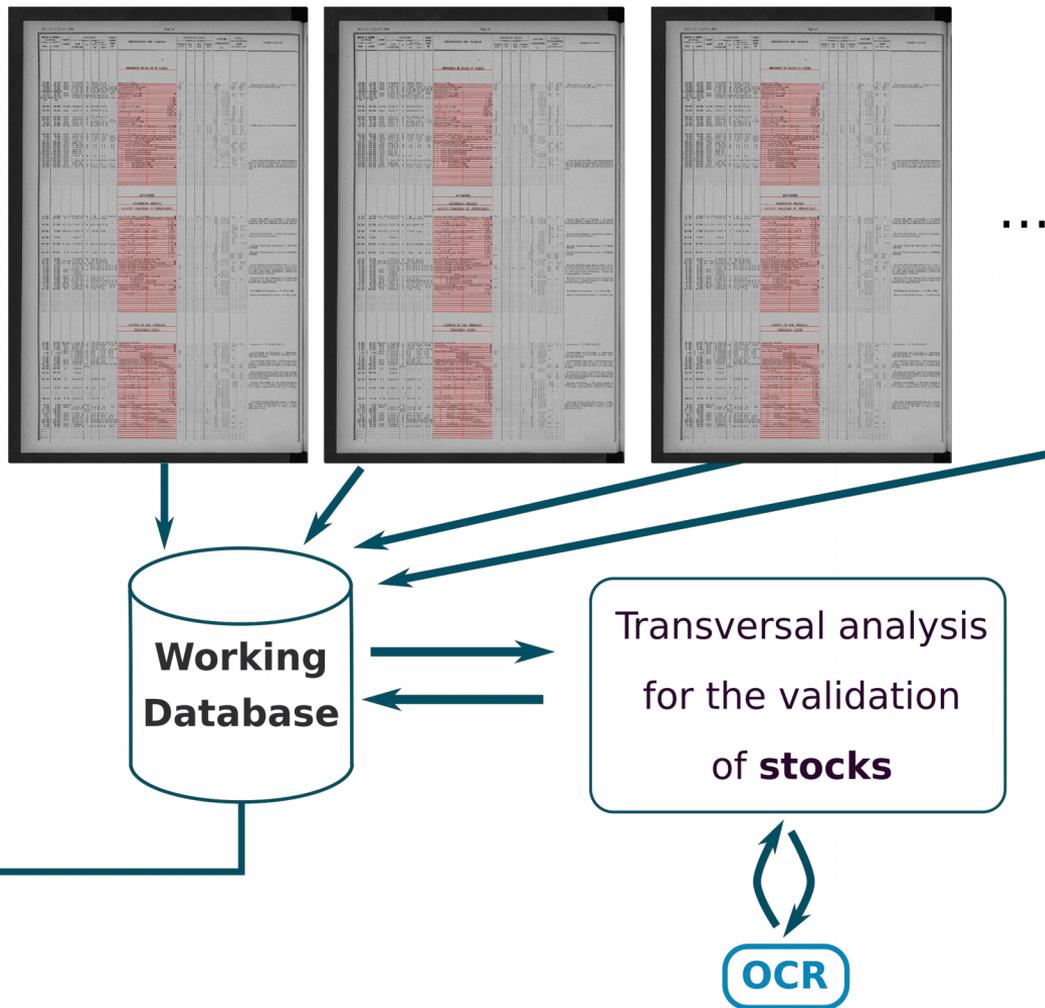
First implementation

Conclusion

Iterative process

Structural Analysis

- Columns (reliable)
- Sections (reliable)
- Stocks



Introduction

Documents characteristics

Collection recognition strategy

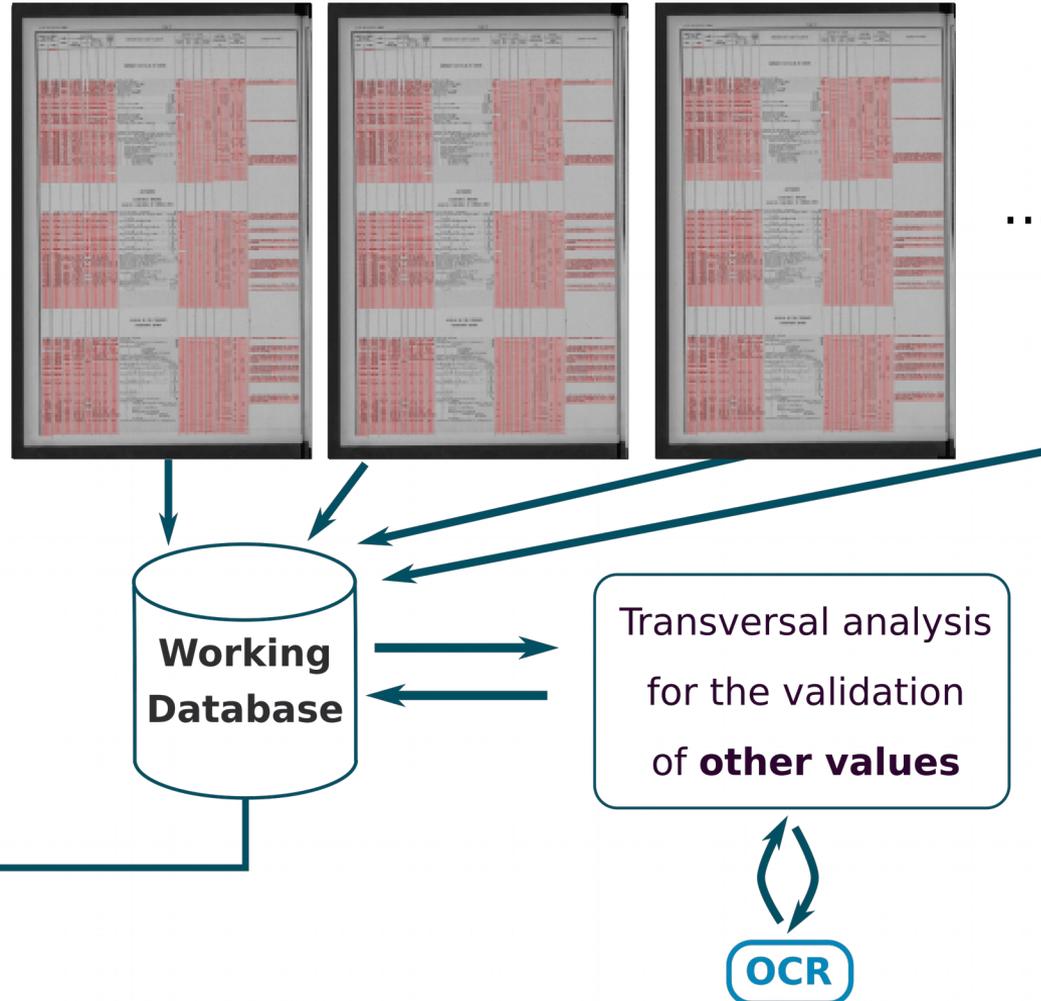
First implementation

Conclusion

Iterative process

Structural Analysis

- Columns (reliable)
- Sections (reliable)
- Stocks (reliable)
- Other values



Introduction

Documents characteristics

Collection recognition strategy

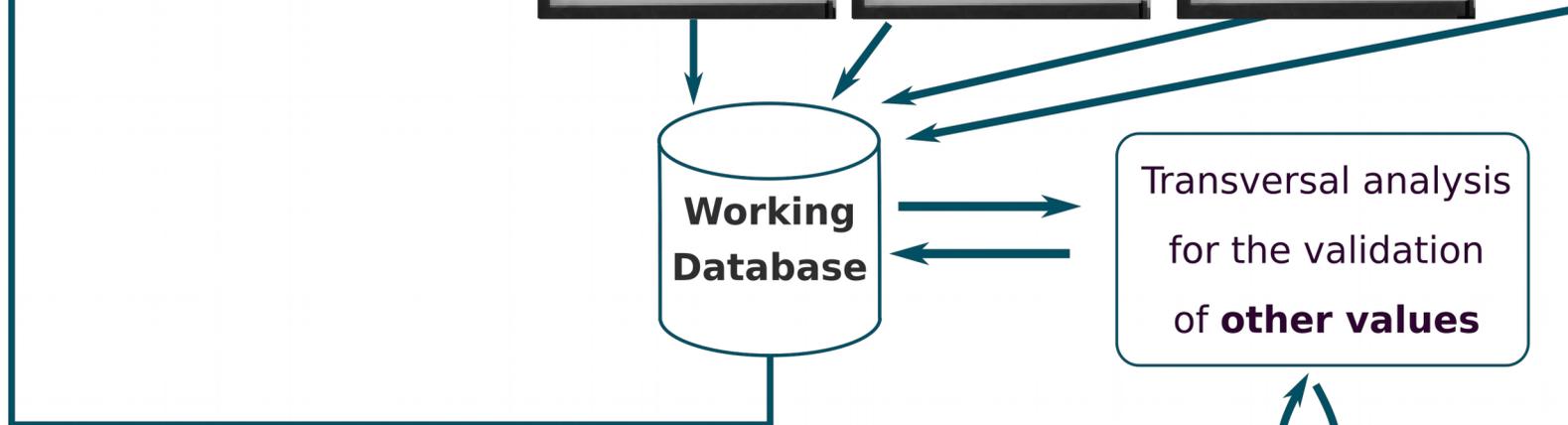
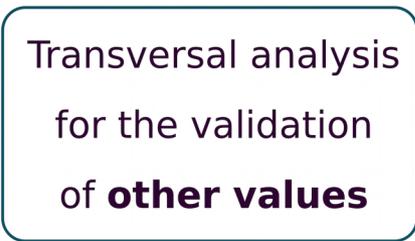
First implementation

Conclusion

Iterative process

Structural Analysis

- Columns (reliable)
- Sections (reliable)
- Stocks (reliable)
- Other values (reliable)



Introduction

Documents characteristics

Collection recognition strategy

First implementation

Conclusion

First implementation

-

Analysis of individual pages

Introduction

Documents
characteristics

Collection
recognition
strategy

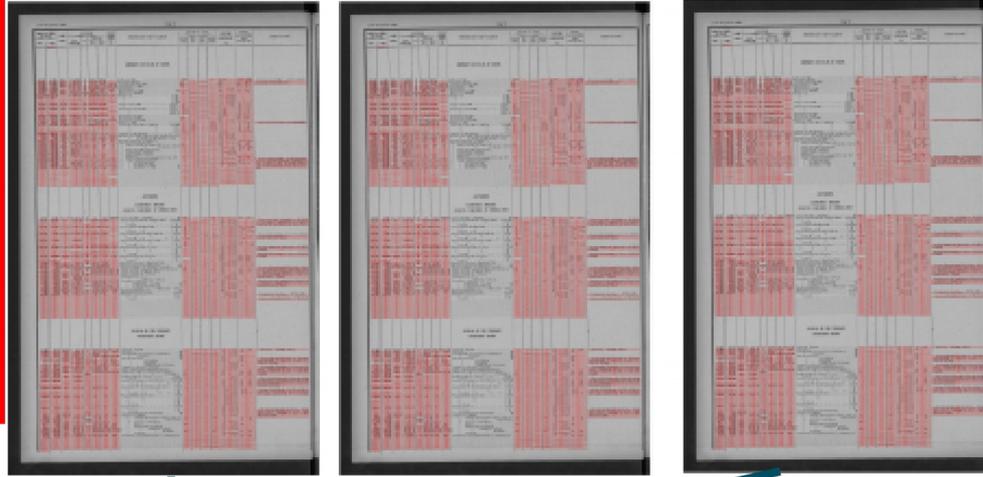
**First
implementation**

Conclusion

Intégration dans la stratégie globale

Analyse structurale

- Colonnes fiabilisées
- Sections fiabilisées
- Titres fiabilisés
- Autres champs fiabilisés



Analyse transversale pour la fiabilisation des **autres champs**



Introduction

Documents characteristics

Collection recognition strategy

First implementation

Conclusion

Purpose and Principle

Purpose:

- First step of the global strategy
- Data generation for OCR training
- Study of the vocabulary contained in the documents

Principle: Combination of a deep learning and a syntactical approach

- Text-lines extraction (deep learning part)
- Structural description of the documents (syntactical part)

Introduction

Documents
characteristics

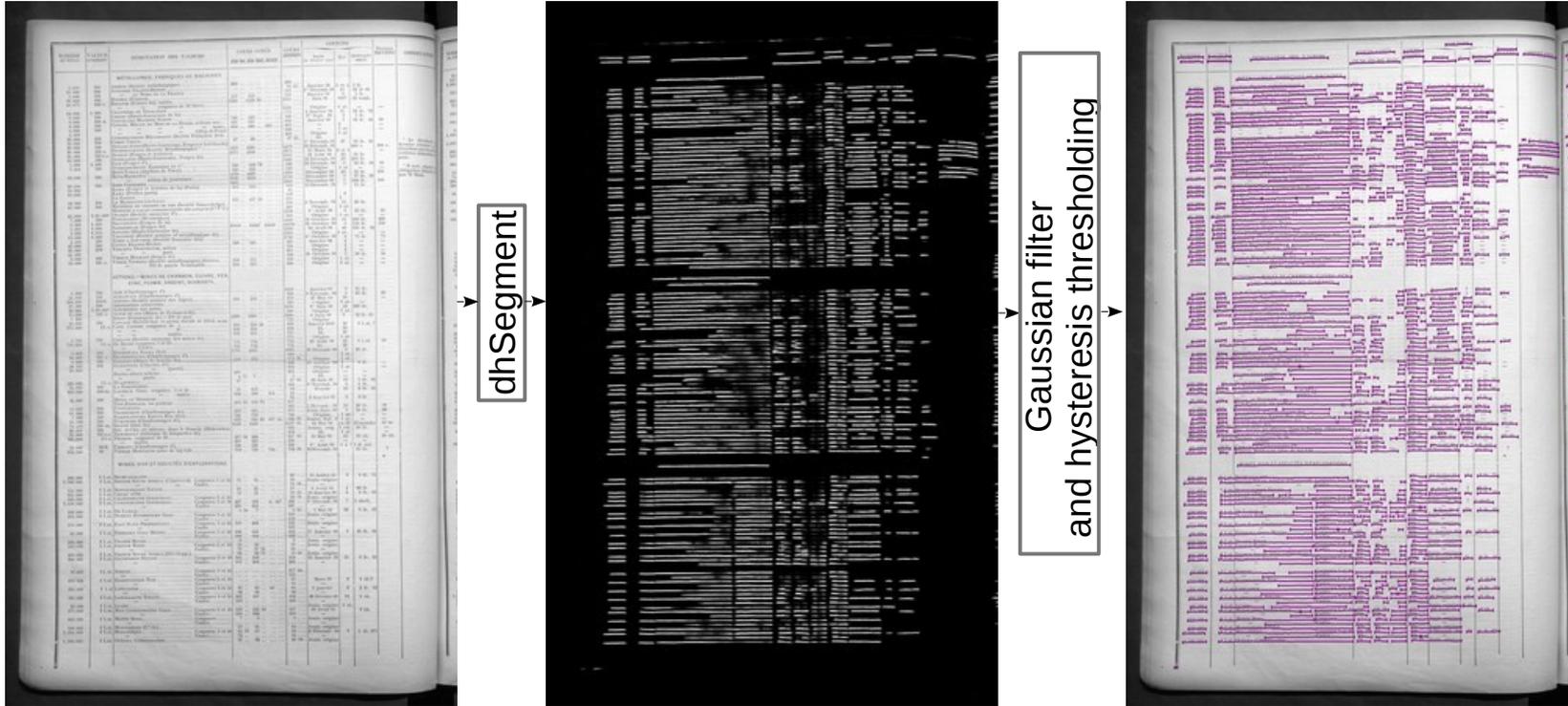
Collection
recognition
strategy

**First
implementation**

Conclusion

Deep Learning for text-lines extraction

dhSegment [S. A. Oliveira (2018)]



- Prediction of the classes of each pixel thanks to a CNN
- Simple post-processing for the extraction of text-lines

Introduction

Documents characteristics

Collection recognition strategy

First implementation

Conclusion

Grammatical Description

structure ::=

banniereDeTitre **Banniere** **EnsColonnes** &&
lesDonnees **EnsLignes**.

...

Introduction

Documents
characteristics

Collection
recognition
strategy

First
implementation

Conclusion

Grammatical Description

structure ::=

*banniereDeTitre Banniere EnsColonnes &&
lesDonnees EnsLignes.*

lesDonnees ::=

lesSections EnsSection.

lesSections [Section | R] ::=

*uneSection Section &&
lesSections R.*

uneSection Section ::=

*intituleDeSection UnIntitule &&
titreBoursier EnsTitreBoursier&&
consSection UnIntitule EnsTitreBoursier Section.*

...

MARDI 5 JANVIER 1924

MONTANT EN RUBRES DE TITRES		VALER		COUPONS				Etendue précédent	DÉSIGNATION DES VALEURS	COURS DU JOUR				CLOTURE	RELEVÉ	OBSERVATIONS
EMIS	ADMIS	nominal	DATE de dernier payé	N°	MONTANT NET	BRUT	base	Premier cours		Plus bas	Plus haut	Dernier cours	PRÉCÉDENT (A)	du 1 ^{er} Janv. 1923		
EMPRUNTS DE VILLES ET DIVERS																
55.000	55.000	500 fr.	1 ^{er} Août 1920	15	Fund.	Fund.	RUSSIA 5 % 1912	32	90	51 24	85 25 92 50	Belaruss 4 % 1902 — Coupons n° 42 détachés à la Côte le 3 Décembre 1923.
350.000	350.000	400 fr.	1 ^{er} Mai 14	6	8 fr. 40	8 fr. 40	RUSSIA 4 % 1911	100 50	103	97	
55.000	55.000	500 fr.	25 Nov. 23	12	1 75 %	1 75 %	CORNICI 3 1/2 % 1902	1600	1900	1627	
55.000	55.000	500 fr.	1 ^{er} Oct. 23	52	5 fr. 405	5 fr. 40	ARMÉNIEN 3 1/2 % 1897	840	825	825	
1400000	1.400000	500 fr.	1 ^{er} Oct. 23	85	0 lire 50	1 lire 50	FRANCOIS 3 %	31	24 25	
15.000	15.000	500 fr.	1 ^{er} Dec. 23	42	10 fr.	10 fr.	HESSINGEN 4 % 1900	381	51 23	378	
175.152	175.152	500 fr.	14 Sept. 17	85	1 lire 50	1 lire 50	MOSCOW 5 % 1886-1900	72	72	
80.853.310	80.853.310	lire	1 ^{er} Janv. 24	85	1 lire 50	1 lire 50	NAPLES 5 %	10	24 11 23	
.....	c. 100	10	24 11 23	
.....	c. 300	10	24 11 23	
100.000	100.000	Rb. 100	14 Sept. 17	29	2 R. 870	2 R. 25	ODessa 4 1/2 % 1903	40	24 11 23	
.....	c. 1000	40	24 11 23	
160.000	160.000	500 fr.	15 Oct. 17	31	10 7.6871	11 fr. 25	PETROGRAD 4 1/2 % 1902	50	52	
.....	comp. 5	50	52	
.....	comp. 10	50	52	

Introduction

Documents characteristics

Collection recognition strategy

First implementation

Conclusion

Grammatical Description

structure ::=

*banniereDeTitre Banniere EnsColonnes &&
lesDonnees EnsLignes.*

lesDonnees ::=

lesSections EnsSection.

lesSections [Section | R] ::=

uneSection Section &&

lesSections R.

uneSection Section ::=

intituleDeSection UnIntitule &&

titresBoursiers EnsTitreBoursier&&

consSection UnIntitule EnsTitreBoursier Section.

...

The image shows a page from a financial publication, likely a stock market listing, dated April 4, 1928. The page is titled "COTATIONS" and "COURSES DES JOURS". It contains a dense grid of data, organized into several sections. The top section is "EMPRETS DE BILLES ET JIVERS", followed by "ACTIONS", "ASSURANCES, BANQUES (SOCIÉTÉS FONCIÈRES ET IMMOBILIÈRES)", and "CHEMINS DE FER, TRANSPORTS (TRANSPORTS DIVERS)". Each section contains multiple columns of data, including company names, prices, and other financial metrics. The text is in French and is printed in a small, dense font. The page is numbered "Page 3" in the top right corner.

Introduction

Documents
characteristics

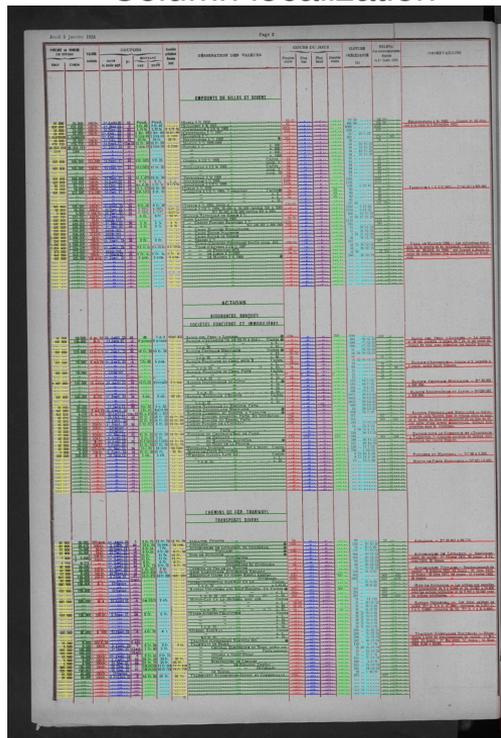
Collection
recognition
strategy

First
implementation

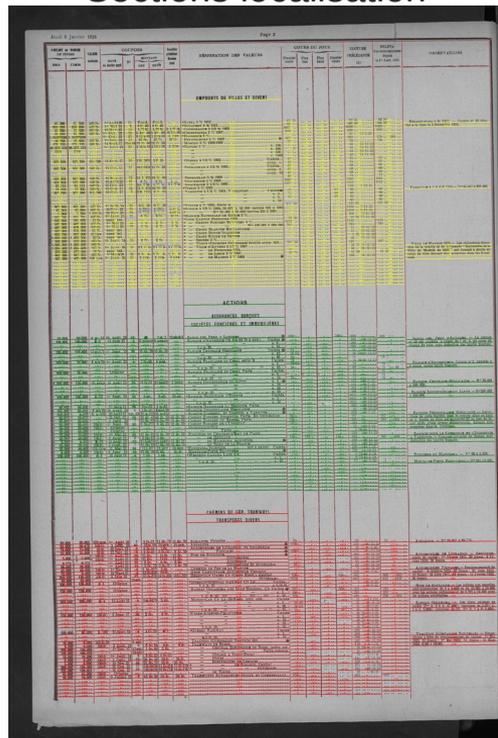
Conclusion

Structural Recognition

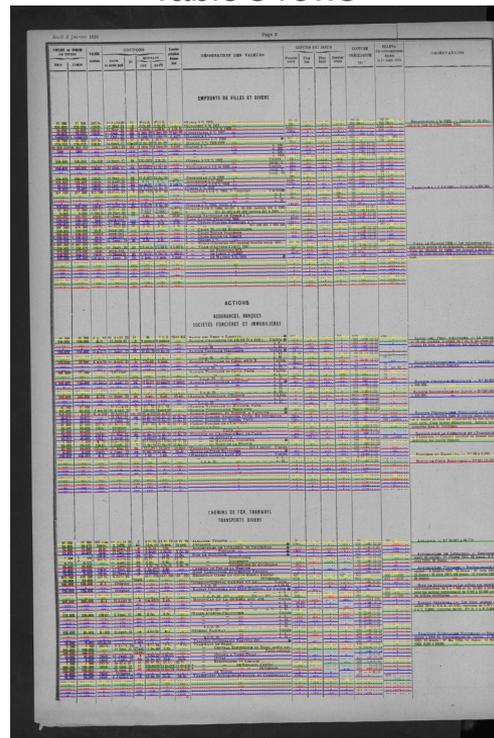
Column localization

A scanned document page with a table structure. The table has multiple columns and rows. The columns are highlighted with vertical colored lines (red, green, blue, yellow) to indicate localization. The table contains text in French, including headers like 'PROJET DE BILAN' and 'REVENUS DE BILAN ET DIVERSES'.

Sections localisation

A scanned document page with a table structure. The table has multiple columns and rows. The rows are highlighted with horizontal colored lines (yellow, green, red) to indicate section localization. The table contains text in French, including headers like 'PROJET DE BILAN' and 'REVENUS DE BILAN ET DIVERSES'.

Table's rows

A scanned document page with a table structure. The table has multiple columns and rows. The rows are highlighted with horizontal colored lines (yellow, green, red) to indicate row localization. The table contains text in French, including headers like 'PROJET DE BILAN' and 'REVENUS DE BILAN ET DIVERSES'.

Result's format :
(a text-line)

```
{ "type" : "line",  
  "polyline" : "1032,5175;1039,5175;1052,5173;1066,5173; ...",  
  "column" : "2",  
  "sectionTitle" : "False",  
  "sectionNumber" : "2",  
  "rowNumber" : "15",  
  "title" : "False" }
```

Introduction

Documents characteristics

Collection recognition strategy

First implementation

Conclusion

First evaluation : Column Recognition

Introduction

Documents characteristics

Collection recognition strategy

First implementation

Conclusion

Période	Du 09/01 au 14/01 1899 (page 2)	Du 14/01 au 31/03 1899 (page 2)	Du 02/01 au 15/04 1924 (page 2)	Du 03/01 au 31/03 1939 (page 2)
Nb d'images	6	64	72	64
Nb d'images avec le bon nb de colonnes détectées	5	58	71	58
Nb d'images avec un nb erroné de colonnes détectées	1	6	1	6

Notes :

- Purpose: - having a first idea of the system performances
- identifying the errors
- Errors will be corrected by the global strategy
- A more precise metric is necessary to correctly assess the performances of our system

Transversal analysis to improve columns recognition

02/01/1924

A scanned document page from 02/01/1924, featuring a table with 16 columns. The table is filled with dense, small text, likely representing a ledger or record book. The columns are clearly defined by vertical lines.

16 columns

03/01/1924

A scanned document page from 03/01/1924, featuring a table with 15 columns. The table is filled with dense, small text, likely representing a ledger or record book. The columns are clearly defined by vertical lines.

15 columns

04/01/1924

A scanned document page from 04/01/1924, featuring a table with 16 columns. The table is filled with dense, small text, likely representing a ledger or record book. The columns are clearly defined by vertical lines.

16 columns

- Error detection :
 - column number
 - OCR results for the columns names
- Identifying the missing column
- Modifying the structural recognition

Introduction

Documents characteristics

Collection recognition strategy

First implementation

Conclusion

Conclusion

Context

- Stock exchange daily list
- A considerable number of documents

Difficulties

- Some documents are strongly deteriorated
- Complex logical structure
- Stock identification

Strategy

- Iterative process
- Traversal analysis to improve the recognition
- A Strategy adaptable to other collection

First implementation: Analysis of individual pages

Introduction

Documents
characteristics

Collection
recognition
strategy

First
implementation

Conclusion

OCR 1 Specifications

Input

...	
TURBIE-SUR-MER, Compagnie du Littoral.	20/01/1899
TURBIE-SUR-MER, Compagnie du Littoral.	21/01/1899
TERRAINS D'ARLES.....	23/01/1899
TERRAINS D'ARLES.....	24/01/1899
TERRAINS D'ARLES.....	25/01/1899
...	
TAVERNES POUSSET ET ROYALE RÉUNIES.	30/01/1899
TAVERNES POUSSET ET ROYALE RÉUNIES.	

49th stock
of the section
«VALEURS DIVERS»
for all the documents

Output

range 1:

- **Interval:** 01/01/1899 - 21/01/1899
- **Transcription Hypothesis:** “Turbie-sur-Mer, Compagnie du Littoral”

range 2 :

- **Interval:** 23/01/1899 – 29/01/1899
- **Transcription Hypothesis:** “Terrains d’Arles”

range 3 :

- **Interval:** 30/01/1899 - ...
- **Transcription Hypothesis:** “Tavernes Pousset et Royale Réunies”

...

Using OCR 1

01/01/1899	46	RAFFINERIE SAY	act. de jouissance
	47	SUD-RUSSE, Soude.....	
	48	TAVERNES POUSSET ET ROYALE RÉUNIES.....	
	49	TURBIE-SUR-MER, Compagnie du Littoral.....	
		...	
21/01/1899	46	RAFFINERIE SAY	
	47	SUD-RUSSE, Soude.....	
	48	TAVERNES POUSSET ET ROYALE RÉUNIES.....	
	49	TURBIE-SUR-MER, Compagnie du Littoral.....	
		...	
23/01/1899	46	RAFFINERIE SAY	act. de jouissance
	47	SUD-RUSSE, Soude.....	
	48	TAVERNES POUSSET ET ROYALE RÉUNIES.....	
	49	TERRAINS D'ARLES.....	
	50	TURBIE-SUR-MER, Compagnie du Littoral.....	
		...	
28/01/1899	46	RAFFINERIE SAY.....	
	47	SUD-RUSSE, Soude.....	
	48	TAVERNES POUSSET ET ROYALE RÉUNIES.....	
	49	TERRAINS D'ARLES.....	
	50	TURBIE-SUR-MER, Compagnie du Littoral.....	

Transversal analysis

OCR 1: put together same stocks

Detecting potential errors:

- 1 stock appear just for 1 day
- 1 stock is missing just for 1 day
- 1 stock disappear and reappear during a short period

Processing:

- Ask for a confidence score (OCR 2)
- Take a decision during the next iteration
- Ask an human operator

Collection characteristics

1) Rules for a same day

NOMBRE de titres	VALEUR nominale	DÉSIGNATION DES VALEURS	COURS COTÉS		COURS précédents	COUPONS		
			plus bas, plus haut, dernier			DATE du dernier payé	Nos	MONTANT BRUT
		MÉTALLURGIE, FABRIQUES DE MACHINES						
5.000	500	ARIÈGE (Société métallurgique) ..	245 .. 253	230				
35.000	100	ATELIERS FRANCO-RUSSES ..	88 .. 89	88	Janvier 98	3 et 4	5 fr.	

- **Interdependence between values**

ex : 'Cours plus bas' \leq 'Cours plus haut'

- **Take the date into account**

ex : 'Date du dernier payé' \leq 'Date of the day'

Collection characteristics

2) Rules between days

03/01/1939

Mardi 3 Janvier 1939 Page 2

MONTANT ou NOMBRE DE TITRES		VALEUR nominale	COURS de compens ^{on}	COURS de REPORTS	DÉSIGNATION DES VALEURS	ÉCHÉANCES des PRIMES	COURS DU JOUR				CLOTURE précédente	COUPONS	
ÉTÉS	ALTES						Premier cours	Plus bas	Plus haut	Dernier cours		DATE DU DERNIER PAYSÉ	MONTANT NET
410.000	410.000	100 fr.	174 ..	1 20	Equateur (Compagnie Générale de l'), ex-c. 11, N° 1 à 410.000.....X P. D. d/20f	fin Dec. d/20f	173 50	174 ..	174 ..	172 ..	170 ..	12 Octobre 1938	9 464 Ⓜ

04/01/1939

Mercredi 4 Janvier 1939 Page 2

MONTANT ou NOMBRE DE TITRES		VALEUR nominale	COURS de compens ^{on}	COURS de REPORTS	DÉSIGNATION DES VALEURS	ÉCHÉANCES des PRIMES	COURS DU JOUR				CLOTURE précédente	COUPONS	
ÉTÉS	ALTES						Premier cours	Plus bas	Plus haut	Dernier cours		DATE DU DERNIER PAYSÉ	MONTANT NET
410.000	410.000	100 fr.	174 ..	1 20	Equateur (Compagnie Générale de l'), ex-c. 11, N° 1 à 410.000.....X P. D. d/20f	f. c. d/20f	174 50	174 ..	174 ..	172 ..	170 ..	12 Octobre 1938	9 464 Ⓜ

- **Interdependence between values**

ex : 'Dernière Clôture' is equal to 'Dernier COURS' of previous day

- **Stable Values**

ex: 'VALEUR NOMINALE'

- **Change at a precise date**

ex : 'COURS de Compensation'