

Projets

Mise en situation professionnelle

4 Topics → 5 sujets

- + Classification (2 sujets)
- + Analyse de documents (1 sujet)
- + Fouille de processus (1 sujet)
- + Big Data (1 sujet)

Travail demandé

+ Rapport (20 pages maximum sans les références)

- Etat de l'art :
travaux existants, étude comparative, synthèse ...
- Méthodologie :
Votre méthode, Code, screen shots ...

+ Soutenance (~20 minutes de présentation + 10 minutes de questions)

- Partie théorique
- Partie pratique
- Préparer bien votre discours

Evaluation

	Rapport	Soutenance
Qualité de présentation	3 points	5 points
Qualité de rédaction	5 points	--
Réponses aux questions	--	3 points
Code	4 points	--

Sujet 1

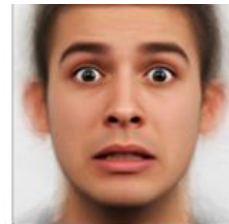
Détection d'émotions

+ Objectif

- Assigner l'émotion convenable à un commentaire



Sadness



Fear



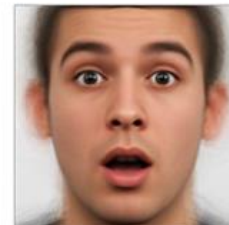
Anger



Neutral



Happiness



Surprise



Disgust

Sujet 1

Détection d'émotions

+ Dataset

- 6000 commentaires
- 3 labels (anger, fear, joy)
- Format : csv

Lien : <https://www.kaggle.com/datasets/abdallahwagih/emotion-dataset>

+ Méthodologie

- Explorer et appliquer des méthodes de machine learning (SVM, Naïve Bayes, arbre de décision...)
- Diviser les données en ensembles de train et de test
- Evaluation des résultats et étude comparative
- Analyse de résultats

Sujet 2

Classification de texte (fake/real news)

+ Objectif

- Détecter la classe associée (**fake** ou **real**) à un texte

+ Dataset

- 21k real data
- 23k fake data
- Format : csv
- Data → title, text, subject, date

Lien : <https://www.kaggle.com/datasets/clmentbisailon/fake-and-real-news-dataset>

Sujet 2

Classification de texte (fake/real news)

+ Méthodologie

- Analyse statistique de données (sujets/ date)
- Pré-traitement de données (stop words, encodage...)
- Application de méthodes de machine learning
- Division les données en ensembles de train et de test
- Evaluation des résultats et étude comparative
- Analyse de résultats

Sujet 3

Séparation de flux documentaires

+ Objectif

- Retrouver la structure initiale d'un flux de documents

+ Dataset

- 9 articles scientifiques fusionnés en 1 pdf
- 81 pages
- 1 fichier csv pour la structure

Lien : <https://ao.univ-lr.fr/index.php/s/gRiAFar47iPzR4b>

Sujet 3

Séparation de flux documentaires

+ Méthodologie

- Considérer le flux comme une séquence de paires de pages
- Extraire le texte à l'aide de l'OCR
- Identifier des descripteurs de continuité/rupture
- Appliquer un vote majoritaire

↗ Descripteurs de continuité

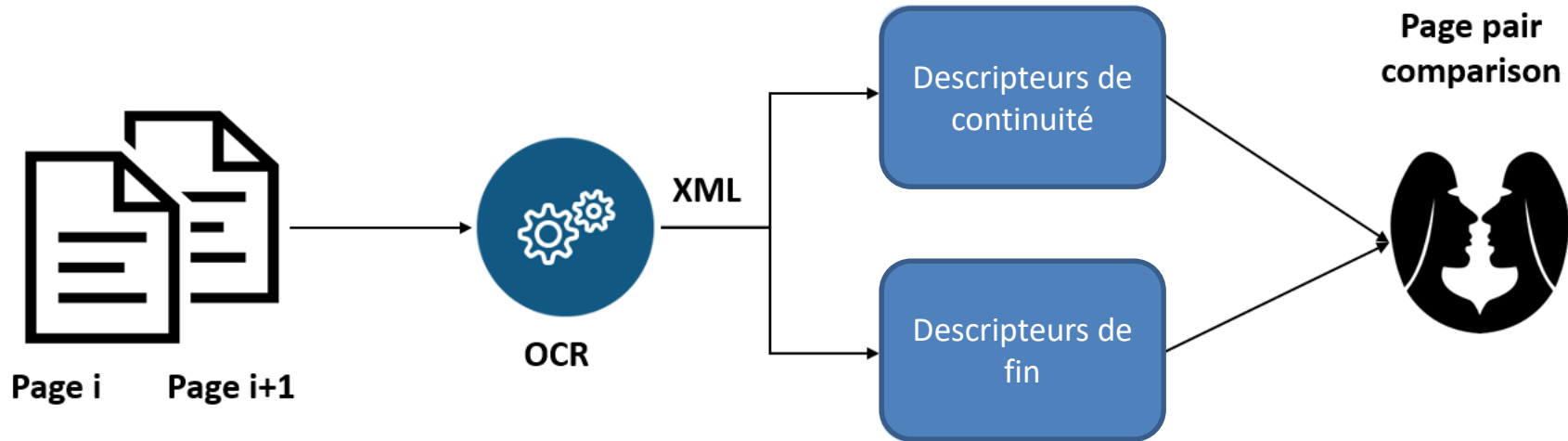
- > Pagination
- > Entités nommées

↗ Descripteurs de rupture

- > Signe de début : *adresses mail, abstract...*
- > Signe de fin : *bibliographie*

Sujet 3

Séparation de flux documentaires



+ Résultats attendus

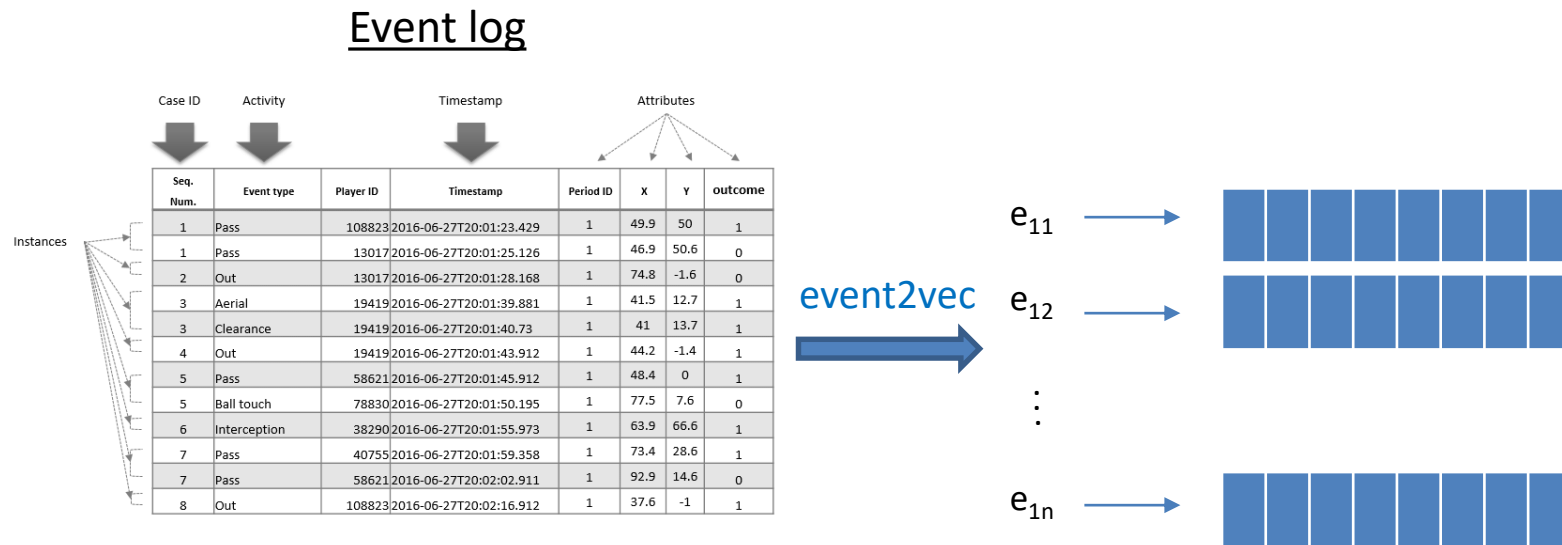
- Fichier csv (séquence de pages – séquence de valeurs booléennes)
- Évaluation : précision (documents, pages), erreurs (segmentation, fusion)

Sujet 4

activity2vec

+ Objectif

- Convertir des activités en vecteurs à l'aide d'un modèle word2vec que vous devez entraîner sur quelques millions de traces



Sujet 4

activity2vec

+ Travail demandé

- Construction d'un « multidomain » event log
- Étude statistique (domaine, activités, variants)
- Apprendre un modèle word2vec sur les activités
- Calculer la similarité entre les activités

+ Links to event logs

Domaine	Liens de téléchargement
Education	MoocData
	Open Learning Analytics OU Analyse Knowledge Media Institute The Open University
	SDP@US1.0 - Dataset for Student Dropout Prediction (figshare.com)
	https://github.com/riiid/ednet
	Educational Process Mining (EPM): A Learning Analytics Data Set - UCI Machine Learning Repository
Système de recommandation	https://www.kaggle.com/datasets/retailrocket/ecommerce-dataset?select=events.csv
Librairie numérique	https://ao.univ-lr.fr/index.php/s/3tYqfkg4RrB7JrN
Aachen université	https://www.pads.rwth-aachen.de/cms/PADS/Forschung/Event-Logs/~bcbswu/-Un-Fair-Event-Logs/

Sujet 5

Big Data

PART -I- vers un big data responsable

+ Impact environnemental du big data

- 2% de la consommation de l'électricité mondiale
- Grandes émissions de consommation électrique
 - 80g de CO₂ par kwh en France
 - 460g de CO₂ par kwh en Allemagne

+ Objectifs

- Détailler l'impact du big data sur l'environnement
- Proposer des solutions envisageables (matérielles et logicielles)

Sujet 5

Big Data

PART -II- SKACK is the new SMACK

+ Spark Mesos Akka Cassandra Kafka → S Kubernetes ACK

- Dans ce projet on se propose de construire une étude théorique pour éprouver l'intérêt de l'architecture SKACK vs SMACK.

+ Objectifs

- Décrire l'architecture SKACK
- Décrire l'architecture SMACK
- Etude comparative

Répartition des étudiants

Sujet	Etudiants
1	
2	
3	
4	
5	

Deadline et ordre de passage

Rapports - date limite de soumission : mercredi 08/11/2023 à 13h00

Présentation - ordre de passage : vendredi 10/11/2023

Sujet	Heure de passage
1	13h15
2	13h45
3	14h15
4	14h45
5	15h15



D'ici, on voit + loin !



univ-larochelle.fr