# Early Detection of Psychotic Disorders: the Role of Emotions

**Summer School on Digital Humanities**
**September 13TH TO September 14TH, UNAM, Mexico**
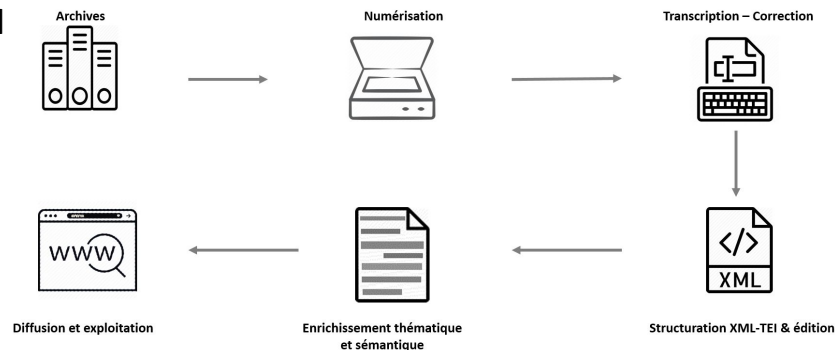
Motasem Alrahabi, ObTIC

# Presentation

- 2010 : PhD in Linguistic Engineering from Sorbonne University, Paris.
- 2010 - 2018 : ICTE Lecturer - Sorbonne University UAE.
- 2018 - today: Research Engineer in Digital Humanities - ObTIC, Sorbonne Univ., Paris.
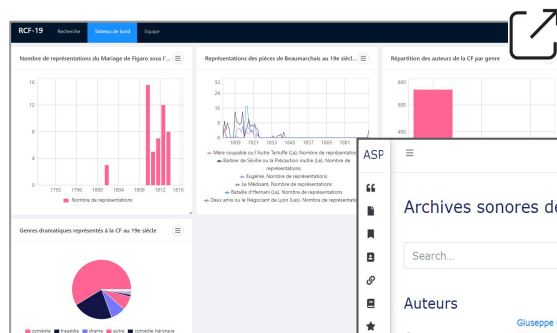  - AI, NLP, Semantic and Discursive Analysis, Digital Publishing, ICT for Education.

# L'Observatoire des textes, des idées et des corpus (ObTIC)

- **ObTIC** (former LabEx **OBVIL**) is a project team dedicated to Digital Humanities at Sorbonne University.
- Collaboration with the **SCAI** (Sorbonne Center for Artificial Intelligence) and the **Datalab** (National Library of France).
- ObTIC draws on the expertise acquired in:
  - Production and digital edition of data (see **OBVIL Library**).
  - Design and experimentation of text mining tools (TAL, AI, corpus analysis, textometry, etc.) for the human and social

# Digital Publishing and Databases

- Automatic File Conversion Tool (Teinte)
- OBVIL Digital Library
- Valentin Haüy Digital Library (AVH)
- Sound heritage of poetry (ASP)
- Registers of the French Comedy (RCF-19)
- Revolutionary Opéra-comique Database (OCD)

# Text Exploration and Mining

- Obvie: Corpus Linguistic Analysis
- Elicom: Explore correspondence and letters
- Ariane: Semantic analysis of texts
- Tanagra: Mapping place names in texts
- Summarizer: Summarizing scientific articles
- Pandore: The toolbox for digital humanities

# APAISE Project

# Context

- Project [APAISE](): *Apprentissage Profond pour l'Analyse Informatisée de la Subjectivité et des Émotions dans les troubles psychotiques émergents*.
    - Deep Learning for Computational Analysis of Subjectivity and Emotions in Emerging Psychotic Disorders
- Funding: Fondation de France (Psychiatric Disease Research [program]())
    - Start: September 2023
- Involved teams:
    - ObTIC - Sorbonne University: Motasem Alrahabi, Jean Marie Tshimula (soon)
    - INSERM: Marie-Odile Krebs, Julien Descles, Valeria Lucarini
    - Brest University Hospital Centre: Michel Walter, Christophe Lemey, Deok-Hee Kim-Dufor
- Goals:
    - Early detection of psychotic disorders: the role of emotions (and the retelling of memories).
    - Better management of patients at risk and slowing down their evolution towards chronicity.

# State of the Art

- There are numerous NLP studies that analyse *language* in the emerging psychotics disorders.
- Identify discriminating anomalies to predict the evolution of patients with differents methods, on different levels: prosody, syntax, semantics, vocabulary, discourse, dialogue…
- Clinical evolution of psychosis: 3 phases
  - At-risk, Not-at-risk, Psychotic "CAARMS" score [Yung et al., 2005]
- Challenge: reduce the prodromal phase (at-risk) before the chronic phase.

→ [Magaud et al., 2010], [Tanguy et al., 2011], [Register-Brown, Hong LE 2014], [Bedi et al. 2015], [Bazziconi, 2018], [Ratana et al., 2019], [Lucarini et al., 2023]…

# Working Hypothesis

- Clinical studies: the subjective dimension of language expression (emotions, sentiments, perceptions…) could reflect in the patient's speech a disturbed relationship to the world and to themself.
- Few academic works in this field [Tshimula et al., 2022], [Saffar, 2023].
- Hypothesis: the analysis of subjective modalities could play an important role in the early detection of psychosis (during the prodromal phase).

→ Classification problem: given a labeled text as input, what would its class belong to?

| 1 | At-risk (A) |
|---|---|
| 2 | Not-at-risk (N) |
| 3 | Psychotic (P) |
| 4 | Control (C) |

# Data Preprocessing

# Data

- **Our corpus consists of rare psychiatric interviews:**
  - Open dialogues between psychiatrists and patients (15 to 30 years old).
- **About 250 audio interviews (differents patients).**
  - Currently: 134 interviews (≈ 1 million tokens).
- **Patient speech is characterised by:**
  - Verbal pauses and disfluencies, hesitations, silences...
  - Disorganised speech (tangential and incoherent), broken syntax…
  - Low lexical density, short sentences...
  - Particular use of personal pronouns (moi, je, me, mon, ma, mes)
  - Emotions: stress, anger, violence, euphoria, joy, suffering…

| 1 | **At-risk (A)** | **65 texts** |
|---|---|---|
| 2 | **Not-at-risk (N)** | **17 texts** |
| 3 | **Psychotic (P)** | **19 texts** |
| 4 | **Control (C)** | **33 texts** |

# Data Preprocessing

- We conducted a series of preprocessing on the data:
  - Manual transcription of audio files into text format.
  - Anonymisation: identification of named entities with SpaCy, then manual correction.
  - Oral errors are not corrected: unfinished words, morpho-syntax errors, conjugations...
  - Keep verbal disfluencies and hesitations: *ah, euh, hum, hmm, hein, ben, bah, pfff…*

# Data labeling

- Only for patients' speech
- Applied features for each patient (text level analysis):
  - Label #1 → Average of sentence length
  - Label #2 → Average of the personal pronouns
  - Label #3 → Average of the verbal disfluencies
  - Label #4 → Lexical density of vocabulary
  - Label #5 → Subjective modalities: emotions, sentiments, opinions…

# Subjective Modalities

# Data labeling: Subjective Modalities

- Allows us to capture subjective content:
    - Polarity: positive, negative, neutral or mixed.
    - Source and target.
    - Intensity (normal, strong, etc.).
    - Aspects of the analyzed object.

    → [Turney, 2002]; [Wiebe et al., 2005]; [Pang and Lee, 2008]
      [Balahur et al., 2011]; [Zhang and Bing 2017]…



| | | |
|---|---|---|
| Sentence | | |
| **Subjectivity classification :** | Objective | Subjective |
| **Polarity classification :** | | Positive  Negative  Neutral |

- Need for more fine-grained classification:
    - GoEmotions: 27 labels for emotions in English [Demszky et al., 2020].

# Data labeling: Subjective Modalities

- Examples from the dataset (approximate translation):
  - **One time I took my Swiss army knife and went for a walk and I, I just wanted to shove it down my throat.**
    - *Une fois j'avais pris mon couteau suisse et j'étais parti me promener et je, je voulais juste me le planter dans la gorge.* [Violence, Patient 2]
  - **I couldn't tell the difference between, between if I was in a dream or if I was in reality.**
    - *J'arrivais plus à faire la différence entre, entre si j'étais dans un rêve ou si j'étais dans la réalité.* [Hallucination, Patient 15]
  - **I just want to drink, until, finally, [I lose reason], because I have, I am in control of myself all the time.**
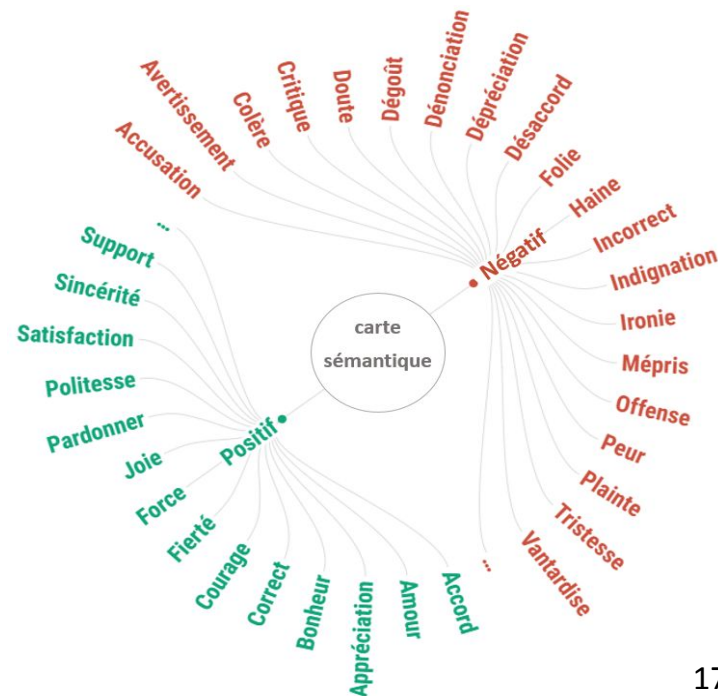    - *J'ai juste envie de boire, jusqu'à, enfin, la déraison, enfin parce que j'ai, je suis tout le temps dans le contrôle de moi-même.* [Addiction, Patient 401]
  - **Each time I thought I had zeros [in exams], I realized that in fact uh they were just kidding me.**
    - *A chaque fois je pensais avoir des zéros, je me rendais compte qu'en fait euh on se foutait juste de ma gueule* [Mockery, Patient 1101]

# Data labeling: Subjective Modalities

- Linguistic Ontology [Alrahabi, 2016, 2021]
  - ≈ 3500 observable markers (patterns)
  - Lexical categories: verbs, adjectives, adverbs, phrases…

- Fine-grained annotations:
  - Classified first as positive, negative or neutral.
  - Grouped into 82 sub-categories: anxiety, stress, anger, violence, joy, suffering…
  - Adaptation to the current project: set up new categories, consider oral speech, existing oral pronunciation errors…

# Data labeling: Subjective Modalities

- Lexicon-based annotation tool ([Textolab](#)):
  - 28770 annotations / 134 texts / 1,050,144 words.

Annotation results can be consulted via a web interface ([Ariane](#))

# Supervised Classification

# Data Representation (Embeddings)

- Vector representation: Camembert LM (https://huggingface.co/camembert-base)
- Embeddings are created using the "CamembertTokenizer" (based on "WordPiece").
- Embeddings are associated with features:
  - 4 linguistic labels
  - 82 emotion labels

    → Dimensions: 768 + 86 = 854

- Embeddings visualized with PCA:
  - No underlying clusters.



Embeddings Visualization (before training)

# Supervised classification

- Use of traditional machine learning algorithms (no enough data for Deep Learning)

**Embeddings**

**Classification**

**Cross-Validation**

# Evaluation of classification models

- LazyPredict library (https://pypi.org/project/lazypredict/)

| # | Model | Accuracy | F1 Score | Time |
|---|-------|----------|----------|------|
| 1 | LGBMClassifier | 0,78 | 0,77 | 1,01 |
| 2 | XGBClassifier | 0,78 | 0,75 | 1,08 |
| 3 | ExtraTreesClassifier | 0,78 | 0,75 | 0,19 |
| 4 | LinearDiscriminantAnalysis | 0,78 | 0,78 | 0,16 |
| 5 | RandomForestClassifier | 0,78 | 0,73 | 0,34 |
| 6 | Perceptron | 0,74 | 0,76 | 0,05 |
| 7 | CalibratedClassifierCV | 0,74 | 0,7 | 0,39 |
| 8 | RidgeClassifierCV | 0,74 | 0,76 | 0,13 |
| 9 | LogisticRegression | 0,74 | 0,76 | 0,17 |
| 10 | SVC | 0,74 | 0,66 | 0,04 |

# Evaluation of the best classification models

- Perform a cross-validation in terms of accuracy:
  - assessing model performance
  - tuning hyperparameters
  - ensuring generalization to new data
  - etc.

| Model | F1 | F2 | F3 | F4 | F5 | Mean Score |
|---|---|---|---|---|---|---|
| LGBMClassifier | 0,78 | 0,74 | **0,81** | 0,70 | 0,69 | 0,75 |
| ExtraTreesClassifier | 0,74 | **0,78** | **0,78** | 0,74 | 0,69 | 0,75 |
| LinearDiscriminantAnalysis | 0,63 | 0,74 | **0,78** | 0,74 | 0,58 | 0,69 |

# Evaluation of the best classification models

- Precision, recall and f-score:
  - evaluate the performance of classification models (quality, completeness and overall performance)
  - offer interpretable measures of a model's performance
  - aid in choosing the most appropriate machine learning algorithm for a given problem
  - etc.

| Model | Precision | Recall | F1-Score |
|---|---|---|---|
| LGBMClassifier | **0,81** | 0,71 | 0,74 |
| ExtraTreesClassifier | 0,77 | 0,71 | 0,73 |
| LinearDiscriminantAnalysis | 0,77 | **0,8** | **0,78** |

# Evaluation of the best classification models

- Confusion Matrix:
  - offers a general overview of the performance of a classification model by summarizing the counts of true positive, true negative, false positive, and false negative predictions



**LGBMClassifier**　　　　**ExtraTreesClassifier**　　　　**LDA**

# Analysing the importance of labels (features)



*Emotions play a very important role*

# Conclusion

# Synthesis: Hybrid Approach

**Data: open dialogues**

Q ........................
A ........................
Q ........................
A ........................
Q ........................

↓

**Preprocessing**

- NER (SpaCy) for anonymization.
- Linguistic labels: sentence length, disfluencies, personal pronouns, lexical density, (later: prosody)
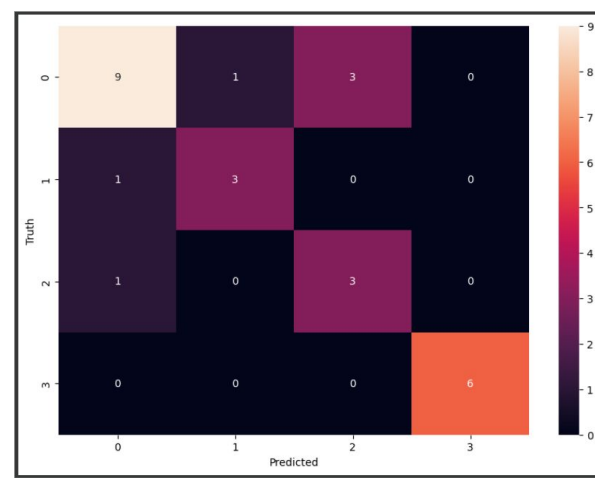- Emotions labels: stress, anger, violence, euphoria, joy, suffering…Lexicon-based tool (Textolab).
- Metadata (later): sex, age, clinical information.
- Embeddings: Camembert Language Model.
- Corpora exploration: Ariane.

↓

**Supervised Classification**

- Supervised Machine Learning (Logistic Regression, Decision Trees, Random Forest, SVM, K-NN…)
- Cross-Validation
- Evaluation: Confusion Matrix, Precision, Recall, F-Score…)
- Features Importance

28

# Preliminary Results and Perspectives

- Recruitment of a postdoctoral fellow (oct. 2023).
- Use data sampling or sliding window techniques.
- Create a multi-label emotion model [Tao et al. 2020], [Demszky et al., 2020].
- Cross with other information:
  - Prosodic analysis: measurement of silence and intonation, in progress with INSERM [Lucarini et al., 2023].
  - Metadata: gender, age, clinical observations (risk of psychosis, consumption of products, etc.).

# Thank you for your attention !

# Bibliographic references

- Alexandre D, Alrahabi M, Gay F., Riguet M. "Le médical et le social: analyse sémantique des rapports de l'immersion d'étudiants de médecine dans le Samu social", in Humanités Numériques Littéraires, sous la dir. de Didier Alexandre, Paris, Classiques Garnier, 2021
- Alrahabi M, Ariane: dispositif de fouille et de lecture synthétique de textes, Actes de DigitAl Humanities and cuLtural herItAge: data and knowledge management and analysis, Jan 2021, Montpellier, France.
- Bedi G, Carrillo F, Cecchi GA, Slezak DF, Sigman M, Mota NB, Ribeiro S, Javitt DC, Copelli M, Corcoran CM. Automated analysis of free speech predicts psychosis onset in high-risk youths. NPJ Schizophr. 2015.
- Bazziconi P-F, 2018, thèse de doctorat de médecine DOI: 10.13140/RG.2.2.16573.82400, Advisor: Christophe Lemey; Michel Walter
- Bazziconi P-F, Bleton L, et al.. L'Information Psychiatrique, 2019/2 (95), 89-94.
- Demszky, Dorottya & Movshovitz-Attias, Dana & Ko, Jeongwoo & Cowen, Alan & Nemade, Gaurav & Ravi, Sujith. (2020). GoEmotions: A Dataset of Fine-Grained Emotions. 4040-4054. 10.18653/v1/2020.acl-main.372.
- Lucarini et al. 2021. "Conversational metrics, psychopathological dimensions and self-disturbances in patients with schizophrenia" European Archives of Psychiatry and Clinical Neuroscience.
- Register-Brown K, Hong LE. Schizophrenia research. 2014 Dec; 160(1-3):20-6.
- Tao, J., Fang, X. Toward multi-label sentiment analysis: a transfer learning based approach. J Big Data 7, 1 (2020).
- Yung AR, Yuen HP, McGorry PD, Phillips LJ, Kelly D, Dell'Olio M, Francey SM, Cosgrave EM, Killackey E, Stanford C, Godfrey K, Buckby J. Mapping the onset of psychosis: the Comprehensive Assessment of At-Risk Mental States. Aust N Z J Psychiatry. 2005