

Guide pour soumettre une réponse

Partie 1 du challenge IA Santé
CentraleSupélec & ILLUIN Technology
12/1/2022 - 27/1/2022

[Introduction](#)

[Description du challenge IA santé](#)

[Description de la partie 1 du challenge IA : Développer une solution d'IA capable de détecter des termes médicaux et les relations entre eux dans des rapports médicaux.](#)

[Dataset](#)

[Concept médical](#)

[Les assertions](#)

[Les relations](#)

[Calcul du score](#)

[Comment soumettre son code pour être évalué ?](#)

[Qu'est-ce-que devra contenir l'output des prédictions à envoyer à votre coach ?](#)

[Qu'est-ce-que vous devrez faire à la fin du datathon ?](#)

Introduction

Description du challenge IA santé

Ce challenge IA santé est composé de 4 parties :

- Développer une solution d'IA capable de détecter des termes médicaux et les relations entre eux dans des rapports médicaux.
- Concevoir une solution capable de choisir les filtres pertinents pour la construction de cohortes de patients, présents dans les rapports médicaux.
- Concevoir les maquettes d'un démonstrateur pour interagir avec les résultats.
- [Optionnel] Développer le démonstrateur

Ce document a pour objectif d'expliquer aux étudiants comment soumettre leurs résultats pour la 1ère partie du challenge IA afin qu'ils soient évalués dessus.

Description de la partie 1 du challenge IA : Développer une solution d'IA capable de détecter des termes médicaux et les relations entre eux dans des rapports médicaux.

L'objectif de cette partie est de développer une solution d'IA capable de :

- Extraire des **concepts médicaux** à partir de rapports de patients.
- Classifier les problèmes médicaux (une typologie de concept médical) extraits selon leur typologie **d'assertion** : présent, absent, probable, etc.
- Classifier les **relations** qui lient les concepts médicaux entre eux.

Dataset

Un **dataset d'entraînement**, annoté, constitué de 170 rapports médicaux sera partagé au début du challenge.

Un **dataset de validation**, constitué de 128 rapports médicaux sera partagé au début du challenge, mais sans les annotations. Cela permettra aux équipes de s'évaluer sur le dataset de validation tout au long du challenge.

Un **dataset de test**, constitué de 128 rapports médicaux, ne sera pas partagé aux étudiants et permettra d'évaluer les équipes à la fin du challenge.

Concept médical

Les concepts médicaux à détecter sont les suivants :

- **Problème médical** : des phrases qui contiennent des observations faites par des patients ou des cliniciens, sur le corps ou l'esprit du patient qui sont considérées comme anormales ou causées par une maladie.
- **Traitements** : des phrases qui décrivent les procédures, les interventions et les substances données à un patient dans le but de résoudre un problème médical.
- **Tests** : des phrases qui décrivent des procédures, des panels et des mesures qui sont faites à un patient ou à un fluide corporel ou un échantillon afin de découvrir, d'exclure ou de trouver plus d'informations sur un problème médical.

Les assertions

Les assertions sont un attribut des concepts de problèmes médicaux qui sont détectés dans la tâche d'extraction de concepts. En tant que telle, cette tâche consiste à classer chaque problème médical dans une catégorie d'assertions. Chaque problème médical sera attribué à l'une des six catégories d'assertions. La tâche de classification ne consiste pas à déterminer si le patient a eu le problème, mais à déterminer ce que la note affirme être le problème médical en fonction du contexte dans lequel elle est utilisée.

Les 6 catégories d'assertion sont les suivantes :

- Présent
- Absent
- Possible
- Conditionnel
- Hypothétique
- Non associé à un patient

Les relations

La tâche de relation consiste à déterminer le type de relation qui existe entre deux concepts dans le texte (le cas échéant). Les relations s'appuient sur les concepts de problème médical, de traitement et de test qui ont déjà été marqués. La tâche consiste à identifier comment les problèmes médicaux sont liés aux traitements, aux tests et aux autres problèmes médicaux dans le texte.

Il y a 3 types de relation :

- Problème médical → Traitement
- Problème médical → Test
- Problème médical → Problème médical

Calcul du score

Le score final sera la moyenne des 3 f1-scores macro associés à chacune des 3 tâches (détection des concepts, des assertions et des relations).

Par exemple si vous obtenez les 3 f1 scores macro suivants :

- 0.86 pour les concepts
- 0.71 pour les assertions
- 0.75 pour les relations

Alors votre score final sera de $(0.86 + 0.71 + 0.75)/3$, soit 0.77.

Comment soumettre son code pour être évalué ?

Explications globales :

1. Faire tourner votre solution d'IA (regroupant les 3 tâches) sur les rapports médicaux du dataset de validation (pendant le challenge) ou le dataset de test (à la fin du challenge).
2. Transmettre l'output de vos prédictions à votre coach ILLUIN pour qu'il puisse évaluer vos résultats

Qu'est-ce-que devra contenir l'output des prédictions à envoyer à votre coach ?

Contenu

C'est exactement le même format que les données dans le dataset d'entraînement.

Nous allons quand même le détailler. L'output doit contenir 3 dossiers à la racine nommés de la manière suivante :

- **concept** : ce dossier contient les prédictions sur la tâche de **détection des concepts**. Il contient autant de fichiers que de rapports médicaux (soit 477 pour le dataset de test). Un fichier dans ce dossier doit respecter la convention de nommage suivantes : `<nom_du_rapport_medical>.con`
- **ast** : ce dossier contient les prédictions sur la tâche de **classification des assertions**. Il contient autant de fichiers que de rapports médicaux (soit 477 pour le dataset de test). Un fichier dans ce dossier doit respecter la convention de nommage suivantes : `<nom_du_rapport_medical>.ast`
- **rel** : ce dossier contient les prédictions sur la tâche de **classification des relations**. Il contient autant de fichiers que de rapports médicaux (soit 477 pour le dataset de test). Un fichier dans ce dossier doit respecter la convention de nommage suivantes : `<nom_du_rapport_medical>.rel`

Le format des fichiers de prédiction (`<nom_du_rapport_medical>.con` ; `<nom_du_rapport_medical>.ast` ; `<nom_du_rapport_medical>.rel`) doit respecter exactement le même que ceux présent dans le dataset d'entraînement.

Qu'est-ce-que vous devrez faire à la fin du datathon ?

1. Envoyer vos prédictions à votre coach technique pour obtenir votre score final.
2. Envoyer votre code à votre coach technique, via un lien github privé. Ceci lui permettra d'évaluer l'implémentation de votre approche technique.