

SOI 2020- DataByter

DataByter è una piattaforma per la gestione di dataset per Machine Learning, sviluppata per il progetto di Sistemi Orientati ad Internet del A.A. 2020/2021. Di seguito la discussione e relativa implementazione dei requisiti funzionali del progetto.

Requisiti funzionali

Accesso

I punti 1 e 2 riguardano l'accesso al sito e la possibilità di registrarsi. Per accedere al sito è richiesto un nome utente e una password. È possibile anche registrarsi, inserendo una mail, uno username e una password. Non è possibile registrarsi utilizzando una mail o uno username già registrato. E però possibile aggiornare la propria password. Quest'ultima operazione richiede però che l'utente sia a conoscenza sia del proprio username e e-mail inseriti durante la procedura di registrazione. In caso contrario non sarà possibile aggiornare la propria password. Sia nella fase di registrazione che di aggiornamento password è necessario che l'utente inserisca due volte la stessa password per evitare possibili errori di battitura. Tutti i controlli vengono effettuati lato server. I dati di accesso degli utenti vengono salvati in chiaro all'interno del database. Questo rappresenta un problema per quanto riguarda la sicurezza, che potrebbe essere risolto tramite l'utilizzo di algoritmi di crittazione. Un altro punto critico non gestito è il salvataggio dello username all'interno del sessionStorage. Un utente malevolo infatti potrebbe modificare il valore salvato e agire come un altro utente

Progetto

I punti dal 3 al 8 riguardano la struttura dei progetti. L'utente può creare un nuovo progetto inserendo un nome, una descrizione e un "obiettivo" di dimensione del dataset. Quest'ultimo indica quale dovrebbe essere la dimensione minima del dataset per essere ritenuto adatto all'addestramento di un modello. In ogni caso questo è solo un valore indicativo e non produce nessuna limitazione al progetto. Si possono selezionare due tipi di progetto: Immagine o Testo. I progetti Immagine prevedono l'inserimento di due o più campi (fields) di cui uno dovrà essere selezionato come il campo dei labels, cioè il campo su cui i modelli dovranno addestrarsi. Tutti i campi a parte quello dei labels sono campi Immagine, potranno quindi contenere solo immagini. I progetti Testo invece prevedono anche la possibilità da parte dell'utente di definire un tipo per il campo: Testo, Numerico, Binario, Data. In questo modo gli utenti che vorranno aggiungere o modificare delle entry del progetto dovranno attenersi a una struttura più rigida, permettendo una più facile manutenibilità del dataset e un processo di Data Cleaning più rapido. In entrambi i tipi di progetto occorre inoltre inserire tutti i possibili valori del campo Label. In questo modo è il creatore del progetto a gestirne la struttura e a decidere se il dataset potrà essere multi-label o meno. Una volta creato il progetto non può essere modificato ma solo eliminato.

Entry

I punti 9, 10 e 11 riguardano la gestione delle entry degli utenti nei vari progetti. Ogni utente ha la possibilità di visualizzare o eliminare qualsiasi progetto, anche non creato da lui. La pagina del progetto contiene una tabella ordinata con i campi e i valori aggiunti dagli utenti, oltre che a dei grafici che mostrano lo stato del progetto rispetto alla propria soglia minima e a bilanciamento (balance) tra i label delle entry. È spesso preferibile infatti avere dei dataset il più possibile bilanciati per evitare episodi di overfitting durante l'addestramento. Ogni utente può anche aggiungere, modificare o eliminare un qualsiasi campo. La pagina di aggiunta di un nuovo campo contiene una form strutturata in base ai tipi di campi aggiunti durante la creazione del progetto. Nel caso di modifica di una entry la pagina è simile alla precedente con la differenza che i campi di input vengono popolati di default con il valore salvato nella entry. L'utente può visualizzare anche lo storico di una specifica riga, vedendo ad esempio le modifiche che sono state fatte nel tempo da ciascun utente. Infine è possibile anche scaricare il dataset sotto forma di file CSV: nel caso di progetti Testo i dati sono quelli effettivi, mentre per i progetti Immagine in ogni cella viene salvata una stringa Base64 che rappresenta l'immagine salvata.

Requisiti tecnici

Database

Si è scelto di utilizzare MongoDB, uno dei più noti database non relazionali, per il salvataggio dei progetti e delle relative entry. La notevole libertà offerta da questo tipo di database ben si sposa con la natura dinamica di questi progetti: risulta infatti molto semplice la creazione di progetti con strutture diverse. Un altro aspetto molto utile è il fatto che le entry vengono salvate come document, i quali presentano lo stesso formato di un oggetto JS. È quindi possibile eseguire delle query, ottenere il risultato e inviarlo direttamente al client senza aver bisogno di particolari modifiche alla struttura. Inoltre MongoDB, come molti altri database, offre la possibilità di eseguire query asincrone, permettendo una gestione più efficiente delle chiamate.