

# DSD+



## Maintenance Report

Corso di Ingegneria, Gestione ed Evoluzione del Software

**Docente:**

Andrea De Lucia

**Partecipanti:**

Daniele Fabiano 0522501738

**Tutor:**

Francesco Paolo D'Antuono 0522501767

Gilberto Recupito

## Indice

|  |           |
|--|-----------|
| <b>Descrizione del problema</b>            | <b>4</b>  |
| <b>Sistema Attuale</b>                     | <b>5</b>  |
| Descrizione del sistema                    | 5         |
| Descrizione delle pagine                   | 5         |
| Home page                                  | 5         |
| Customization page                         | 6         |
| Result page                                | 7         |
| Saved results page                         | 7         |
| My profile page                            | 8         |
| Login page & Sign up page                  | 8         |
| Installazione del sistema                  | 9         |
| Installazione locale                       | 9         |
| Installazione containerizzata (Docker)     | 10        |
| <b>Reverse Engineering</b>                 | <b>11</b> |
| System & Object design                     | 11        |
| Package                                    | 11        |
| Package core                               | 11        |
| Package detector                           | 16        |
| Documentazione dei package core e detector | 25        |
| Package expectations                       | 26        |
| Package web_application                    | 28        |
| Gestione dei dati persistenti              | 29        |
| Generazione dello schema dei dati          | 30        |
| Analisi dei Requisiti                      | 31        |
| Requisiti funzionali                       | 31        |
| Requisiti non funzionali                   | 32        |
| Pseudo-requisiti                           | 32        |
| Modello dei casi d'uso                     | 33        |
| UC_01: Upload file                         | 33        |
| UC_01.01: Formato file errato              | 33        |
| UC_01.02: Cambio file                      | 34        |
| UC_02: Data smell detection                | 34        |
| UC_02.01: Parametri non selezionati        | 35        |
| UC_02.02: Cancellazione risultati          | 35        |
| UC_03: Visualizzazione risultati salvati   | 35        |
| UC_03.01: Nessun risultato salvato         | 35        |
| Sequence diagram                           | 36        |
| SD_01: Upload file                         | 36        |
| SD_02: Data smell detection                | 36        |
| SD_03: Visualizzazione risultati salvati   | 37        |
| <b>Change Requests</b>                     | <b>38</b> |
| Change request 1                           | 38        |
| Change request 2                           | 38        |

|   |           |
|---|-----------|
| Change request 3  | 39        |
| <b>Impact Analysis</b>                                  | <b>40</b> |
| Impatto di CR_01  | 40        |
| Grafo diretto per CR_01                                 | 41        |
| Matrice di connettività per CR_01                       | 41        |
| Matrice di raggiungibilità per CR_01                    | 41        |
| Impatto di CR_02  | 42        |
| Impatto di CR_03  | 42        |
| Glossario acronimi per la tracciabilità degli artefatti | 43        |
| Design  | 43        |
| Codice  | 43        |
| Testing   | 45        |
| <b>Selezione Change Request</b>                         | <b>46</b> |
| <b>Forward Engineering</b>                              | <b>47</b> |
| Analisi dei requisiti                                   | 48        |
| Modello dei casi d'uso                                  | 48        |
| UC_CR_01: Upload File                                   | 48        |
| UC_CR_01.01: Formato file errato                        | 49        |
| UC_CR_01.02: Cambio file                                | 49        |
| UC_CR_01.03: Campo form gruppo vuoto                    | 49        |
| UC_CR_03: Visualizzazione risultati salvati             | 50        |
| UC_CR_03.01: Nessun risultato salvato                   | 50        |
| Sequence diagram  | 51        |
| SD_CR_01: Upload file                                   | 51        |
| SD_CR_02: Data Smell Detection                          | 51        |
| SD_CR_03: Visualizzazione risultati salvati             | 52        |
| System & Object Design                                  | 53        |
| Package expectations                                    | 53        |
| Package web_application                                 | 54        |
| Gestione dei dati persistenti                           | 55        |
| <b>Testing</b>  | <b>55</b> |
| <b>Descrizione delle nuove pagine</b>                   | <b>56</b> |
| Home page   | 56        |
| Saved results page                                      | 57        |

# Descrizione del problema

Negli ultimi anni, alcuni studi si sono incentrati sui Data Smell, segnali o indizi sulla presenza di problemi o anomalie latenti nei dati, che richiedono un'attenta e dettagliata analisi. Per questo motivo sono nati alcuni tool che permettono la detection di smell, uno di questi è DSD (Data Smell Detection). Esso si basa sul tool di convalida dei dati open-source Great Expectations e si concentra sul rilevamento degli smell, analizzando un dataset mediante l'utilizzo di tecniche rule-based. DSD+ ha lo scopo di estendere la suite di detection di smell esistente. Inoltre si intende migliorare l'attuale meccanismo di report dei dataset, integrando il calcolo delle metriche di dimensione della qualità dei dati (data quality dimensions). Rispetto al dataset analizzato, sarà possibile visualizzare i risultati dell'esecuzione del tool, con lo scopo di seguire i cambiamenti ottenuti nel tempo a seguito di nuove analisi, per formare una base di indicazioni sul miglioramento/peggioramento di qualità del dataset modificato.

# Sistema Attuale

## Descrizione del sistema

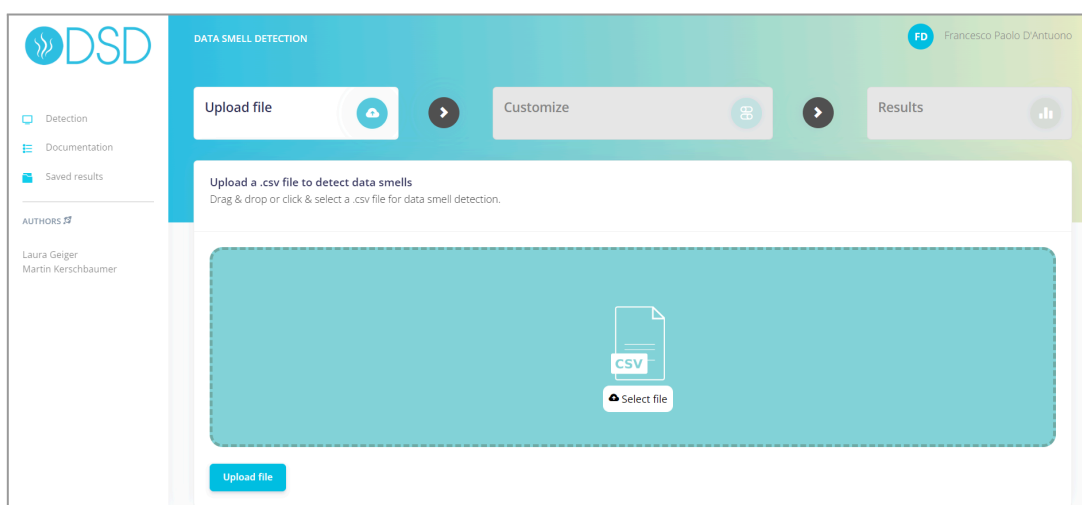
DSD (Data Smell Detection) è un tool sviluppato per consentire l'identificazione di data smell all'interno di dataset. L'implementazione si basa sulla libreria Great Expectations, che è uno strumento open-source progettato per aiutare gli utenti a verificare e documentare la qualità dei dati. Il progetto è stato inizialmente realizzato per la versione python 3.6. Risulta però possibile eseguire il tool sull'ultima versione supportata di python 3.8 (3.8.19). Il tool mette a disposizione una dashboard basata sul template [Argon Dashboard Django](#), che permette di avere una gestione già definita relativamente alla registrazione e autenticazione dell'utente. Inoltre consente di accedere alle funzionalità del sistema che effettuano la detection di data smell sul dataset caricato. Per fornire una maggiore comprensibilità del sistema, sono descritte le principali pagine del sistema.

**Link al repository del progetto:** <https://github.com/CpDant/DSD-plus>

## Descrizione delle pagine

### Home page

L'immagine di seguito mostra come è composta l'home page. A sinistra vi è una barra laterale dove è possibile accedere alle funzionalità principali della web application. In alto a destra vi è la possibilità di accedere al proprio account e al centro viene descritta la funzionalità principale di detection di data smell, con un riquadro che consente l'inserimento del dataset che verrà analizzato.



## Customization page

In questa pagina vengono determinati i parametri specifici che verranno utilizzati dal sistema, per effettuare la detection di data smell in un certo tipo. Inoltre si può scegliere anche di detectare un sottoinsieme di smell specifici, evitando di effettuare analisi inutili per quel tipo di dataset. Le immagini seguenti mostrano ciò che è stato descritto in precedenza.

The screenshot shows the 'Customize your detection' page of the DSD (Data Smell Detection) application. The interface includes a sidebar with navigation links for 'Detection', 'Documentation', and 'Saved results'. The main content area has a top navigation bar with 'Upload file', 'Customize', and 'Results' buttons. Below this, there's a section titled 'Customize your detection' with instructions to choose presets, data smells, and columns. A dropdown menu shows 'Easy Mode' and 'Expert Mode'. Under 'Presettings', there are three buttons: 'Tolerant', 'Medium', and 'Strict'.

This screenshot displays the 'Data smells' configuration section. It is divided into two main categories: 'BELIEVABILITY SMELLS' and 'ENCODING UNDERSTANDABILITY SMELLS'. Under 'BELIEVABILITY SMELLS', there are three checked items: 'Duplicated Value Smell', 'Extreme Value Smell', and 'Suspect Sign Smell'. Each has a 'mostly' parameter set to 0,9 and a 'threshold' parameter set to 3,0. Under 'ENCODING UNDERSTANDABILITY SMELLS', there are three checked items: 'Casing Smell', 'Long Data Value Smell', and 'Missing Value Smell'. Each has a 'mostly' parameter set to 0,9. The 'Casing Smell' also has a 'same\_case\_wordcount\_threshold' set to 2,0. The 'Long Data Value Smell' has a 'length\_threshold' set to 30,0.

This screenshot shows the 'SYNTACTIC UNDERSTANDABILITY SMELLS' configuration section. It includes three checked items: 'Floating Point Number As String Smell', 'Integer As Floating Point Number Smell', and 'Integer As String Smell'. Each has a 'mostly' parameter set to 0,01. The 'Integer As Floating Point Number Smell' also has an 'epsilon' parameter set to 1e-05. Below this, there's a 'Columns' section with seven checked checkboxes: 'FLOAT1', 'FLOAT2', 'INT1', 'INT2', 'STRING1', 'STRING2', and 'STRING3'. A 'Customize' button is located at the bottom of the configuration area.

## Result page

La seguente immagine mostra la pagina dei risultati della detection appena effettuata. La divisione dei risultati viene fatta per le colonne del dataset.

The screenshot shows the 'Results' page of the DSD application. The interface includes a sidebar with navigation links for 'Detection', 'Documentation', and 'Saved results'. The main content area displays the title 'Smell Results By Column Names' and the dataset name 'data\_smell\_testset\_EoWvU7K.csv'. Below this, there are two tabs: 'int2' (selected) and 'string1'. The 'int2' tab shows a table with the following data:

| DATA SMELL TYPE        | TOTAL ELEMENT COUNT | FAULTY ELEMENT COUNT | FAULTY ELEMENT OVERVIEW  |
|------------------------|---------------------|----------------------|--|
| Casing Smell           | 11                  | 6                    | ['abc def ghi', 'abc def ghi', 'casing 1', 'casing 2', 'all lowercase', 'ALL UPPERCASE'] |
| Duplicated Value Smell | 11                  | 2                    | ['abc def ghi', 'abc def ghi']   |

At the bottom of the table, there are two buttons: 'Delete file and result' and 'Upload another file'.

## Saved results page

Questa pagina mostra i risultati salvati nel sistema delle detection effettuate in passato, su dataset differenti. L'immagine seguente mostra come vengono visualizzate le informazioni. Da notare la differenza con la Result page, che genera incoerenza.

The screenshot shows the 'Saved results' page of the DSD application. The interface includes a sidebar with navigation links for 'Detection', 'Documentation', and 'Saved results'. The main content area displays the title 'data\_smell\_testset\_EoWvU7K.csv' and the subtitle 'Selected parameters for detected data smells'. Below this, there are two sections: 'DUPLICATED VALUE SMELL' and 'CASING SMELL'. The 'DUPLICATED VALUE SMELL' section shows the parameter 'mostly: 0.9'. The 'CASING SMELL' section shows the parameters 'mostly: 0.9' and 'same\_case\_wordcount\_threshold: 2.0'. Below these sections, there are two columns with data smells: 'int2' and 'string1'. The 'int2' column shows a table with the following data:

| DATA SMELL TYPE        | TOTAL ELEMENT COUNT | FAULTY ELEMENT COUNT | FAULTY ELEMENT OVERVIEW |
|------------------------|---------------------|----------------------|-------------------------|
| Duplicated Value Smell | 11                  | 2                    | ['B. 8']                |

The 'string1' column shows a table with the following data:

| DATA SMELL TYPE        | TOTAL ELEMENT COUNT | FAULTY ELEMENT COUNT | FAULTY ELEMENT OVERVIEW  |
|------------------------|---------------------|----------------------|--|
| Casing Smell           | 11                  | 6                    | ['abc def ghi', 'abc def ghi', 'casing 1', 'casing 2', 'all lowercase', 'ALL UPPERCASE'] |
| Duplicated Value Smell | 11                  | 2                    | ['abc def ghi', 'abc def ghi']   |

## My profile page

La seguente immagine mostra quali informazioni possono essere visualizzate e modificate dall'utente registrato.

The screenshot shows the 'My account' page for a user named 'CpDant'. The page has a blue header with the DSD logo and a sidebar on the left with navigation links: 'Detection', 'Documentation', 'Saved results', and 'AUTHORS 27'. The main content area is titled 'Hello CpDant' and includes a sub-header 'My account' with the text 'This is your profile page. You can change your username, email address, first name and last name.' Below this, there is a 'USER INFORMATION' section with four input fields: 'Username' (containing 'CpDant'), 'Email address' (containing 'francescopaolo177@gmail.com'), 'First name' (containing 'Francesco Paolo'), and 'Last name' (containing 'D'Antuono'). A 'Save user' button is located at the bottom right of the form.

## Login page & Sign up page

Le seguenti immagini mostrano la pagina di login e di sign up del sistema, per accedere all'account dell'utente registrato e per registrare un nuovo utente.

The screenshot shows the login page of the DSD system. The header features the DSD logo and the title 'Data Smell Detection'. Below the title, there is a welcome message: 'Welcome to data smell detection! On this website you can check your .csv file on data smells. Register and login in to be able to use the entire functions.' The main form is titled 'Add your credentials OR create your own user' and contains three input fields: 'Username', 'Password', and 'Remember me' (a checkbox). A 'Sign in' button is located at the bottom of the form. At the bottom of the page, there are links for 'register' and 'HOME'.

The screenshot shows the sign up page of the DSD system. The header features the DSD logo and the title 'Data Smell Detection'. Below the title, there is a welcome message: 'Welcome to data smell detection! On this website you can check your .csv file on data smells. Register and login in to be able to use the entire functions.' The main form is titled 'Add your credentials or authenticate with an existing account.' and contains five input fields: 'Username', 'First name', 'Last name', 'Email', and 'Password'. A 'Password check' field is also present. A 'Create account' button is located at the bottom of the form. At the bottom of the page, there are links for 'login' and 'HOME'.



# Installazione del sistema

L'installazione del sistema è stata descritta nel README del repository e può avvenire in due modalità diverse:

## Installazione locale

I passi per eseguire l'installazione locale sono:

1. Installare Python3.8 (se già non è installato nel sistema). Se si tratta di sistemi Windows con WSL Ubuntu seguire i seguenti passaggi:

- a. Aggiornare Ubuntu alla versione più recente con i comandi:

```
sudo apt update, sudo apt upgrade
```

- b. Importare su Ubuntu Python 3.8 PPA (una repository che contiene le versioni di Python scaricabili e utilizzabili) con il comando:

```
sudo add-apt-repository ppa:deadsnakes/ppa -y.
```

- c. Aggiornare nuovamente Ubuntu per permettere al sistema di riconoscere la nuova repository aggiunta.

- d. Installare Python 3.8 con il comando:

```
sudo apt install python3.8
```

- e. Infine, bisogna installare alcuni pacchetti aggiuntivi per far funzionare il tutto correttamente con i seguenti comandi:

```
sudo apt install python3.8-dev,  
sudo apt install python3.8-venv,  
sudo apt install python3.8-distutils
```

2. Installare pip con il comando:

```
sudo apt-get install python3.8-pip
```

3. Installare il pacchetto virtualenv, che consente di usufruire del virtual environment di python con il comando:

```
sudo pip install virtualenv
```

4. Avviare il virtual environment con il comando:

```
virtualenv env
```

5. Attivare il virtual environment partendo dalla cartella padre del progetto con il comando:

```
source env/bin/activate
```

6. Dopo il passaggio 5, si passa all'installazione dei requisiti e della libreria di detection di data smell. Innanzitutto spostarsi sulla cartella `data_smell_detection` dalla root del progetto con il comando:

```
cd data_smell_detection
```

7. Installare i requisiti con il comando:

```
pip install -r requirements-dev.in
```

8. (Facoltativo) Per costruire e distribuire i pacchetti di Python utilizzare il seguente comando:

```
python setup.py install
```

NB: Il sistema è configurato per funzionare anche senza l'installazione diretta della libreria, in quanto sono importati i path del progetto della medesima. Nel caso si voglia comunque installare la libreria, ogni volta che si effettua una modifica ad essa è necessario ripetere il passaggio 8.

Per il setup della web application, i passaggi da effettuare sono:

1. Spostarsi sulla cartella *argon-dashboard-django* dalla cartella root del sistema con il comando:

```
cd web_application/argon-dashboard-django/
```

2. Installare i moduli richiesti con il comando:

```
pip install -r requirements.txt
```

3. Per creare le tabelle del database eseguire i seguenti comandi:

```
python manage.py makemigrations,  
python manage.py migrate
```

4. Infine, per avviare la web application eseguire il seguente comando:

```
python manage.py runserver
```

## Installazione containerizzata (Docker)

I passi per riuscire a far funzionare il progetto con Docker sono i seguenti:

1. Assicurarsi di aver [installato](#) Docker
2. Posizionandosi sulla root del progetto, per effettuare la build del medesimo, eseguire il seguente comando: `sudo docker-compose build`
3. Avviare la web application con il comando: `sudo docker-compose up -d`

# Reverse Engineering

## System & Object design

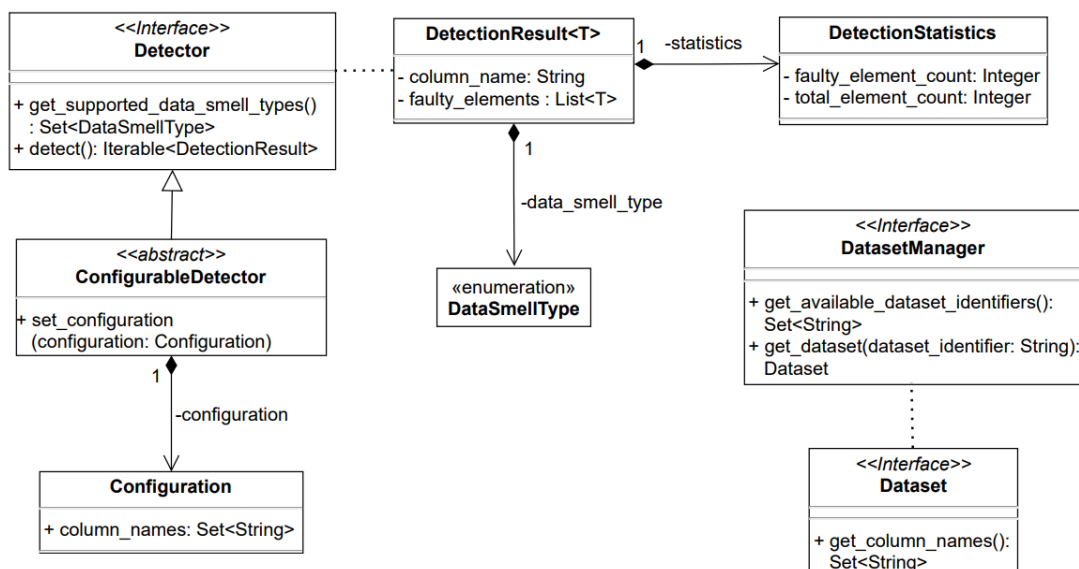
### Package

Il sistema è composto da quattro package differenti, che sono:

- Package **core**: Consente di astrarre il processo di rilevamento dei data smell tramite la classe Detector.
- Package **detectors.great\_expectations**: Contiene la classe GreatExpectationsDetector. Ci riferiremo a questo package, come package detector.
- Package **detectors.great\_expectations.expectations**: Consente di implementare all'interno del sistema i detector specifici dei data smell. Ci riferiremo a questo package, come package expectations.
- Package **web\_application**: Consente di implementare la web application, che utilizza la libreria, per far visualizzare i risultati delle detection su un'interfaccia grafica.

### Package core

Per descrivere il package core è stato ideato un class diagram che riassume tutti gli oggetti e le dipendenze di questa libreria.



Per comprendere al meglio gli oggetti del sistema, di seguito viene riportata una descrizione degli oggetti in tabelle.

| Nome classe      | Descrizione  | Attributi                     |
|------------------|--|-------------------------------|
| Dataset          | Interfaccia che rappresenta il dataset sul quale verrà poi effettuata la data smell detection. | //                            |
| Metodi           |  |                               |
| Nome             | Descrizione  | Parametri e valori di ritorno |
| get_column_names | Consente di prelevare i nomi delle colonne del dataset.  | void -> Set<String>           |

| Nome classe                       | Descrizione  | Attributi                               |
|-----------------------------------|--|---|
| DatasetManager                    | Interfaccia che fornisce metodi per la gestione e la preparazione all'uso del dataset. | //                                      |
| Metodi                            |  |   |
| Nome                              | Descrizione  | Parametri e valori di ritorno           |
| get_available_dataset_identifiers | Consente di prelevare l'insieme di identificatori di dataset.                          | void -> Set<String>                     |
| get_dataset                       | Consente di prelevare il dataset specifico passando un identificatore.                 | dataset_identifier:String<br>-> Dataset |

| Nome classe     | Descrizione   | Attributi   |
|-----------------|---|---|
| DetectionResult | Classe che fornisce un elenco di risultati del processo di detection di data smell. | column_name: String<br>faulty_elements : List<T><br>statistics:<br>DetectionStatistics<br>data_smell_type:<br>DataSmellType |
| Metodi          |   |   |
| Nome            | Descrizione   | Parametri e valori di ritorno   |
| //              | //  | //  |

| Nome classe   | Descrizione  | Attributi   |
|---------------|--|---|
| DataSmellType | Enum che contiene i data smell che i detector possono trovare nei dataset. | Per ogni tipologia di smell è definito un elemento della Enum |
| Metodi        |  |   |
| Nome          | Descrizione  | Parametri e valori di ritorno                                 |
| //            | //   | //  |

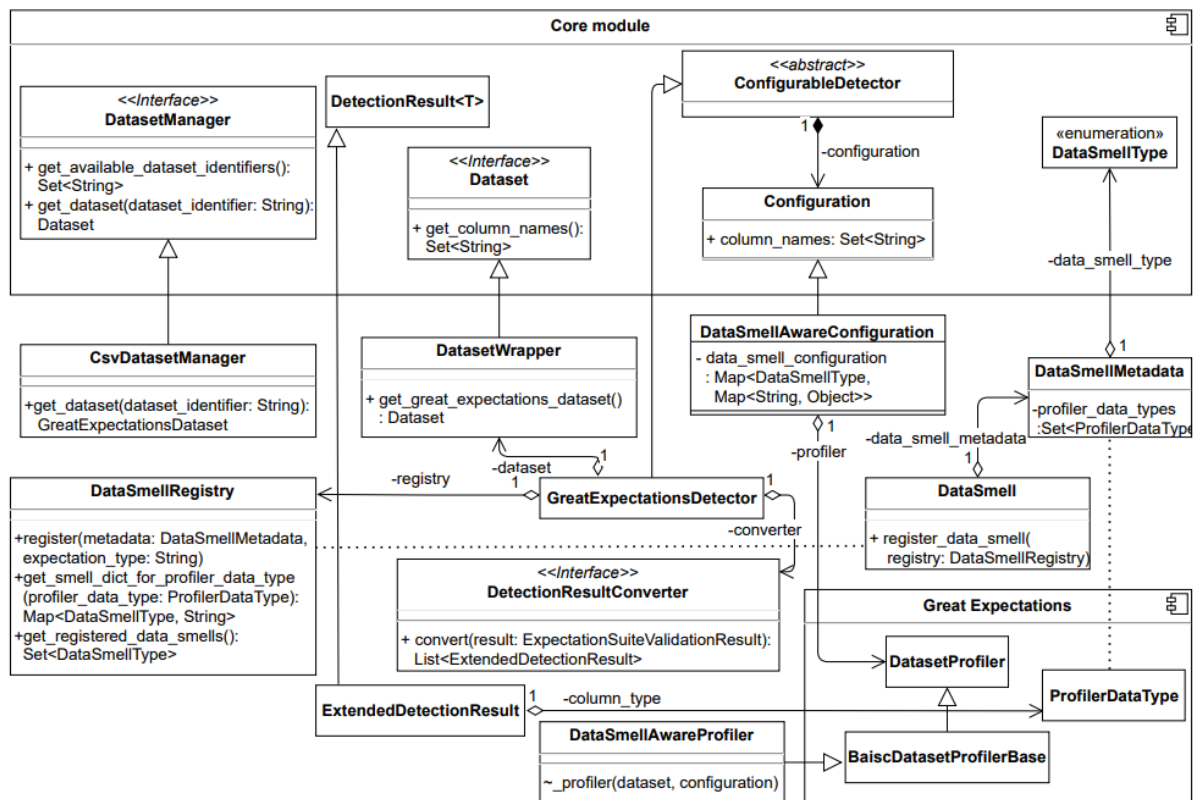
| Nome classe                    | Descrizione   | Attributi                            |
|--------------------------------|---|--------------------------------------|
| Detector                       | Interfaccia che consente di definire i data smell in un particolare dataset.                  |                                      |
| Metodi                         |   |                                      |
| Nome                           | Descrizione   | Parametri e valori di ritorno        |
| get_supported_data_smell_types | Ritorna un insieme di tipi di data smell che il detector può trovare nei dataset.             | void -><br>Set<DataSmellType>        |
| detect                         | Esegue la detection e restituisce i data smell trovati sotto forma di risultati di detection. | void -><br>Iterable<DetectionResult> |

| Nome classe  | Descrizione   | Attributi                     |
|--------------|---|-------------------------------|
| Configurator | Classe astratta che consente di definire i parametri di configurazione per una detection. | column_names :<br>Set[String] |
| Metodi       |   |                               |
| Nome         | Descrizione   | Parametri e valori di ritorno |
| //           | //  | //                            |

| Nome classe          | Descrizione  | Attributi                               |
|----------------------|--|---|
| ConfigurableDetector | Classe astratta che estende la classe Detector per definire un Detector configurabile. | //                                      |
| Metodi               |  |   |
| Nome                 | Descrizione  | Parametri e valori di ritorno           |
| get_configuration    | Restituisce la configurazione usata per il detector.                                   | void -> Configuration                   |
| set_configuration    | Consente di definire una nuova configurazione per il detector.                         | configuration:Configuration<br>n-> void |

## Package detector

Anche qui è stato ideato un class diagram che descrive la composizione degli oggetti del package. Inoltre, viene descritto anche come gli oggetti del package detector estendono gli oggetti del package core e comunicano con alcuni di essi. Da notare che vi è un package esterno, denominato Great Expectations, che è colui che si occupa di accedere agli oggetti utilizzati per svolgere le attività richieste dalle classi a Great Expectations.



Per comprendere al meglio gli oggetti del sistema, di seguito viene riportata una descrizione degli oggetti in tabella (gli oggetti che già sono stati descritti della libreria core non verranno riportati in questa tabella).



| Nome classe                    | Descrizione   | Attributi                     |
|--------------------------------|---|-------------------------------|
| DatasetWrapper                 | Classe che wrappa una istanza di dataset di Great Expectations e estende la classe Dataset. Si trova nel modulo dataset.py. | //                            |
| Metodi                         |   |                               |
| Nome                           | Descrizione   | Parametri e valori di ritorno |
| get_great_expectations_dataset | Restituisce il dataset wrappato.  | void -> Dataset               |

| Nome classe  | Descrizione   | Attributi  |
|--|---|--|
| CsvDatasetManager (nel codice si chiama FileBasedDatasetManager) | Classe che estende la classe DatasetManager e che consente di creare istanze di DatasetWrapper da file CSV. Si trova nel modulo dataset.py. | //   |
| Metodi   |   |  |
| Nome   | Descrizione   | Parametri e valori di ritorno  |
| get_dataset  | Consente di ottenere il dataset importato.  | dataset_identifier:String<br>-><br>GreatExpectationsDataset (anche qui nel codice ritorna in realtà il DatasetWrapper) |

| Nome classe       | Descrizione  | Attributi  |
|-------------------|--|--|
| DataSmellMetadata | Classe che permette di descrivere come effettuare la detection di un data smell specifico. Si trova nel modulo datasmell.py. | data_smell_type:DataSmellType<br>profiler_data_types:Set<ProfilerDataType> |
| Metodi            |  |  |
| Nome              | Descrizione  | Parametri e valori di ritorno  |
| //                | //   | //   |

| Nome classe                               | Descrizione   | Attributi   |
|---|---|---|
| DataSmellRegistry                         | Classe che permette di salvare le istanze di DataSmellMetadata dei data smell implementati. Si trova nel modulo datasmell.py. | //  |
| Metodi                                    |   |   |
| Nome                                      | Descrizione   | Parametri e valori di ritorno   |
| register                                  | Salva un nuovo mapping tra un data smell e la corrispondente aspettativa di Great Expectations.                               | metadata:DataSmellMetadata,<br>expectation_type:String<br>-> //         |
| get_smell_dict_for_profile<br>r_data_type | Permette di prelevare i data smell registrati per uno specifico ProfilerDataType.   | profiler_data_type:Profiler<br>DataType -><br>Map<DataSmellType.String> |
| get_registered_data_smells                | Permette di prelevare un insieme di tipi di data smell registrati.  | -> Set<DataSmellType>   |

| Nome classe         | Descrizione  | Attributi                                       |
|---------------------|--|---|
| DataSmell           | Classe progettata per essere una base per le classi che implementano la detection di data smell. Si trova nel modulo datasmell.py. | data_smell_metadata:Optional<DataSmellMetadata> |
| Metodi              |  |   |
| Nome                | Descrizione  | Parametri e valori di ritorno                   |
| register_data_smell | Semplifica la registrazione delle classi che implementano la detection dei data smell in un DataSmellRegistry.                     | registry:DataSmellRegistry -> //                |

| Nome classe                 | Descrizione  | Attributi  |
|-----------------------------|--|--|
| DataSmellAwareConfiguration | Classe di configurazione che controlla come la detection di data smell viene effettuata dalla classe GreatExpectationsDetector. Inoltre è un'estensione della classe Configuration. Si trova nel modulo detector.py. | data_smell_configuration: Map<DataSmellType, Map<String,Object>> |
| Metodi                      |  |  |
| Nome                        | Descrizione  | Parametri e valori di ritorno                                    |
| //                          | //   | //   |

| Nome classe            | Descrizione  | Attributi                                  |
|------------------------|--|--|
| DataSmellAwareProfiler | <p>Classe che rappresenta un profiler basato su Great Expectations per la detection di data smell.</p> <p>Lo scopo principale è quello di applicare i data smell registrati in una classe DataSmellRegistry, alle colonne del corrispondente DataSmellType. Si trova nel modulo profiler.py.</p> | //   |
| Metodi                 |  |  |
| Nome                   | Descrizione  | Parametri e valori di ritorno              |
| _profile               | <p>Consente di realizzare un ExpectationSuite (descrive i controlli per i data smell, per specifiche colonne in un dataset).</p>   | dataset, configuration -> ExpectationSuite |

| Nome classe             | Descrizione  | Attributi  |
|-------------------------|--|--|
| ExtendedDetectionResult | Estensione specifica di Great Expectations della classe DetectionResult. Si trova nel modulo converter.py. | column_type:ProfilerData Type<br>expectation_kwargs:Dict<str, Any> |
| Metodi                  |  |  |
| Nome                    | Descrizione  | Parametri e valori di ritorno                                      |
| //                      | //   | //   |

| Nome classe              | Descrizione   | Attributi  |
|--------------------------|---|--|
| DetectionResultConverter | Interfaccia che consente di fornire metodi per la conversione da ExpectationSuiteValidationResult a una lista di ExtendedDetectionResult. Si trova nel modulo converter.py. | //   |
| Metodi                   |   |  |
| Nome                     | Descrizione   | Parametri e valori di ritorno  |
| convert                  | Effettua la conversione da ExpectationSuiteValidationResult a istanze di ExtendedDetectionResult.   | result:ExpectationSuiteValidationResult -> List<ExtendedDetectionResult> |

| Nome classe             | Descrizione   | Attributi   |
|-------------------------|---|---|
| StandardResultConverter | Classe che implementa l'interfaccia DetectionResultConverter e dunque la conversione associata. Si trova nel modulo converter.py. | //  |
| Metodi                  |   |   |
| Nome                    | Descrizione   | Parametri e valori di ritorno   |
| convert                 | Implementazione del metodo dell'interfaccia DetectionResultConverter.   | validation_suite_result:ExpectationSuiteValidationResult -> List<ExtendedDetectionResult> |

| Nome classe               | Descrizione   | Attributi                         |
|---------------------------|---|-----------------------------------|
| GreatExpectationsDetector | Sottoclasse della classe ConfigurableDetector. Implementa i metodi definiti sia in Detector che in ConfigurableDetector. Si trova nel modulo detector.py. | //                                |
| Metodi                    |   |                                   |
| Nome                      | Descrizione   | Parametri e valori di ritorno     |
| detect                    | Esegue la detection e restituisce i risultati.  | void -> Iterable<DetectionResult> |



## Documentazione dei package core e detector

La documentazione del codice dei due package può essere generata tramite l'utilizzo del tool [Sphinx](#), attraverso i seguenti passi:

1. Assicurarsi di aver [installato](#) il tool;
2. Posizionarsi nella directory *docs* all'interno di *data\_smell\_detection*;
3. Eseguire il seguente comando per compilare la documentazione:  

```
sphinx-apidoc -o . ../datasmelldetection
```
4. Eseguire il seguente comando per ottenere il formato html:  

```
make html
```
5. Per visualizzare la documentazione del codice, bisogna aprire dal browser il file *index.html*, presente nella directory *\_build/html*

## Package expectations

In questo package sono inserite, per ogni detector implementato, (e dunque per ogni data smell rilevato) una classe specifica che consente di avere una visione completa dei vari detector. Queste classi utilizzano il [pattern di sviluppo di Great Expectations](#), che consente non solo di implementare le expectations personalizzate, ma anche di testarle, utilizzando una struttura già implementata dalla libreria. Inoltre, ad ogni detector è associato un parametro *mostly*, che consente di definire il threshold specifico, per definire che una certa colonna di un dataset contiene lo smell. Le varie classi presenti nel sistema attuale sono:

- **ExpectColumnValuesToNotContainCasingSmell:** Classe che consente di implementare il detector per il Casing Smell. Oltre al parametro *mostly*, è stato inserito un ulteriore parametro, denominato *same\_case\_wordcount\_threshold*, che permette di controllare quante parole devono essere contenute in una stringa, per essere considerate come casing smell.
- **ExpectColumnValuesToNotContainDuplicatedValueSmell:** Classe che consente di implementare il detector per il Duplicated Value Smell. Utilizza la classe ExpectColumnValuesToBeUnique già implementata da Great Expectations.
- **ExpectColumnValuesToNotContainExtremeValueSmell:** Classe che consente di implementare il detector per l'Extreme Value Smell. Oltre al parametro *mostly*, è stato inserito un ulteriore parametro, denominato *threshold*, che permette di capire quando viene identificato un valore come outlier dal detector, fissando una distanza specifica oltre la quale un numero viene considerato un outlier. Il valore di questo parametro deve essere positivo e di default è 3.
- **ExpectColumnValuesToNotContainFloatingPointNumberAsString Smell:** Classe che consente di implementare il detector per il FloatingPointNumberAsStringSmell.
- **ExpectColumnValuesToNotContainIntegerAsFloatingPointNumberSmell:** Classe che consente di implementare il detector per l'Integer As Floating Point Number Smell. Oltre al parametro *mostly*, è stato utilizzato un ulteriore parametro, denominato *epsilon*, che consente di definire il limite per la differenza assoluta tra il floating point e il suo intero più vicino. Il valore di default di questo parametro è 1e-6.
- **ExpectColumnValuesToNotContainIntegerAsStringSmell:** Classe che consente di implementare il detector per l'Integer As String Smell.

- **ExpectColumnValuesToNotContainLongDataValueSmell**: Classe che consente di implementare il detector per il Long Data Value Smell. Oltre al parametro *mostly*, è stato inserito un ulteriore parametro, denominato *length\_threshold*, che permette di determinare il numero di caratteri, per i quali una parola viene considerata “lunga”. Questo parametro deve essere un numero maggiore di 0.
- **ExpectColumnValuesToNotContainMissingValueSmell**: Classe che consente di implementare il detector per il Missing Value Smell.
- **ExpectColumnValuesToNotContainSuspectSignSmell**: Classe che consente di implementare il detector per il Suspect Sign Smell. Oltre al parametro *mostly*, è stato inserito un ulteriore parametro, denominato *percentile\_threshold*, che controlla quali quantili sono calcolati. I quantili sono usati per determinare se la maggioranza dei valori in una colonna è positiva o negativa. Il valore di questo parametro deve essere compreso tra 0 e 1.

Tutti questi detector devono essere inseriti all'interno del DataSmellRegistry, con il metodo *register\_data\_smell()* (descritto nella sezione del package detector). Inoltre, ogni classe avrà *data\_smell\_metadata* come variabile d'istanza, che consente di settare il tipo di data smell e il tipo di dato delle colonne che verranno analizzate dal detector.

## Package web\_application

Questo package ha permesso di implementare il front-end. All'interno di esso sono definite le implementazioni riguardanti sia l'autenticazione, sia la logica applicativa che consente di accedere alle librerie interne del sistema le quali effettuano la detection di data smell. Di seguito è riportata la struttura del package con indicate le componenti principali:

- **web\_application/**
  - **argon-dashboard-django/**
    - **core/** # gestisce i file statici della webapp
      - settings.py # file di configurazione del framework
      - urls.py # definisce gli URL per tutte le pagine
      - static/ # directory che contiene i file statici (css/js)
      - templates/ # template delle pagine da renderizzare
    - **app/** # gestisce la logica applicativa
      - urls.py # definisce gli URL per la logica applicativa
      - views.py # mostra le pagine HTML agli utenti
    - **authentication/** # gestisce la fase di registrazione
      - urls.py # definisce gli URL per la registrazione
      - views.py # mostra le pagine HTML agli utenti
      - forms.py # definisce i form di autenticazione
    - requirements.txt # file di configurazione per le dipendenze
    - manage.py # script di avvio della webapp

A partire da questo package è stato possibile derivare lo schema per la gestione dei dati persistenti come descritto nella sezione successiva.

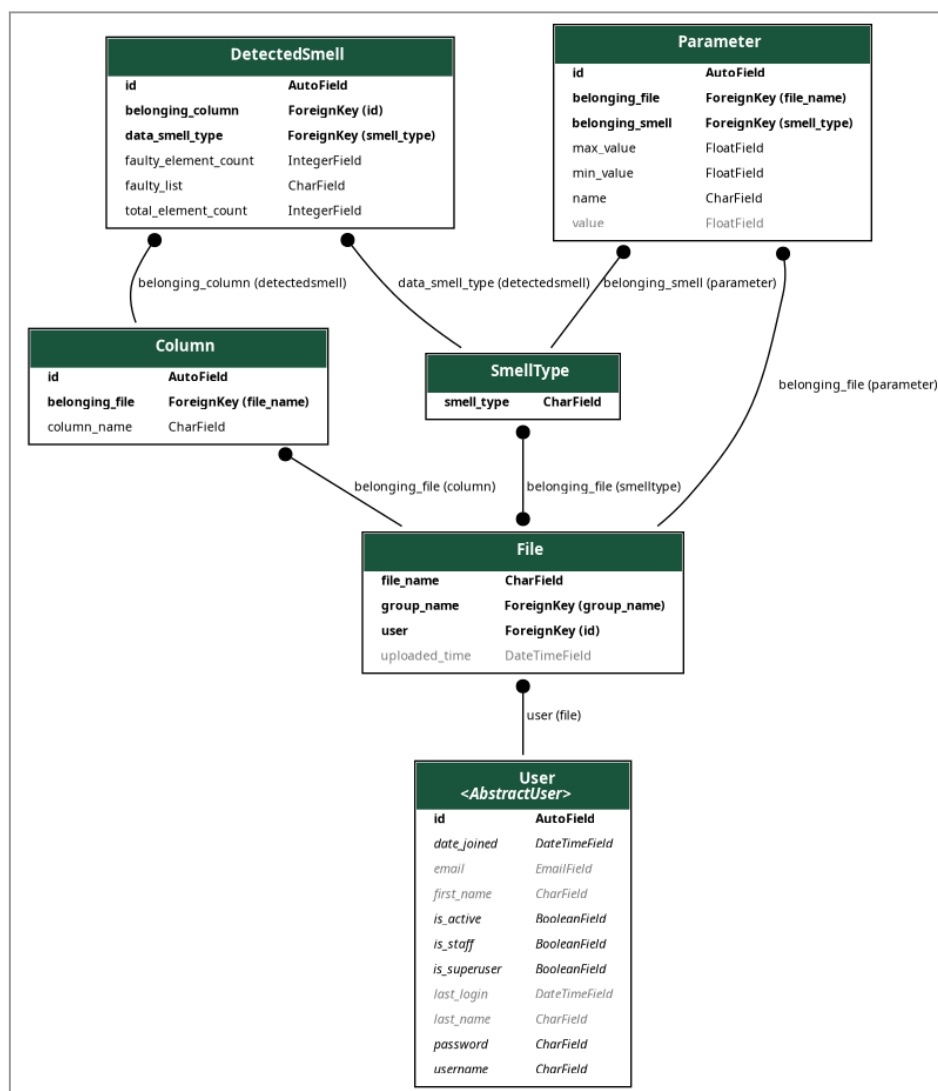
## Gestione dei dati persistenti

Per la gestione dei dati persistenti è stato utilizzato un DBMS relazionale denominato SQLite3. L'utilizzo di esso è stato dovuto da diversi fattori:

- È molto veloce;
- È compatto;
- Ha transazioni ACID;
- Semplice da utilizzare;
- È in grado di interpretare stringhe SQL;

Questo tipo di DBMS viene fornito come libreria da Django, che è il framework utilizzato per lo sviluppo della web application. Infatti, la scelta è stata influenzata dal linguaggio di programmazione utilizzato, cioè Python, e dalla facilità di integrazione di questo tipo di DBMS con esso e il framework della web application.

Nella figura seguente vi è rappresentato lo schema del database.



## Generazione dello schema dei dati

Il framework Django mette a disposizione un insieme di estensioni utili per facilitare la gestione del sistema. Come visto nella precedente sezione è possibile generare lo schema dei dati, grazie all'estensione [Graph models](#), attraverso i seguenti passi:

1. Assicurarsi di aver [installato](#) le django-extensions tramite il comando:

```
pip install django-extensions
```

2. All'interno della sezione *INSTALLED\_APPS* del file *web-application/core/settings.py* inserire la dicitura *django\_extensions*

3. Assicurarsi di aver [installato](#) graphviz

4. Installare pygraphviz con il comando:

```
pip install pygraphviz
```

5. Posizionarsi nella directory dove si trova *manage.py*:

```
cd web_application/argon-dashboard-django
```

6. Per generare lo schema dei dati, eseguire il seguente comando specificando i nomi delle tabelle del database:

```
python manage.py graph_models -a -I DetectedSmell,  
Parameter, Column, SmellType, File, User -o  
schema.png
```

# Analisi dei Requisiti

## Requisiti funzionali

- **RF\_01 - Registrazione al sistema:** L'utente ospite deve poter registrarsi al sistema inserendo le proprie informazioni;
- **RF\_02 - Autenticazione al sistema:** L'utente ospite deve poter autenticarsi con le proprie credenziali;
- **RF\_03 - Modifica informazioni utente:** L'utente registrato deve poter modificare le proprie informazioni;
- **RF\_04 - Logout dal sistema:** L'utente registrato deve poter effettuare il logout dal sistema;
- **RF\_05 - Caricamento del dataset:** L'utente ospite e l'utente registrato devono poter caricare un dataset all'interno del sistema;
- **RF\_06 - Cambio del dataset:** L'utente ospite e l'utente registrato devono poter cambiare il dataset caricato;
- **RF\_07 - Detection di data smell:** L'utente ospite e l'utente registrato devono poter effettuare una detection personalizzata dei data smell all'interno di un dataset appena caricato;
- **RF\_08 - Visualizzazione risultati detection:** L'utente ospite e l'utente registrato devono poter visualizzare i risultati della detection appena effettuata;
- **RF\_09 - Salvataggio risultati detection:** L'utente registrato deve poter salvare i risultati ottenuti dalla detection effettuata per visualizzarli in un secondo momento;
- **RF\_10 - Eliminazione risultati detection:** L'utente registrato deve poter eliminare i risultati ottenuti da una detection, che sono stati precedentemente salvati;
- **RF\_11 - Visualizzazione documentazione smell:** L'utente ospite e l'utente registrato devono poter visualizzare la documentazione dei vari tipi di data smell;

## Requisiti non funzionali

Il sistema esistente ha uno scopo molto specifico e un'utenza ristretta che non corrisponde a quella di un tradizionale cliente, ma a un'utenza interna ad un team di sviluppo. Non conoscendo le specifiche esigenze richieste, non sappiamo quali sono i requisiti non funzionali che devono essere prioritizzati. A seguito della nostra analisi abbiamo comunque individuato alcuni aspetti che potrebbero essere oggetto di migliorie future:

- **Usabilità:** Il sistema spesso non fornisce messaggi chiari all'utente di ciò che sta succedendo nelle fasi di detection e le label di alcuni pulsanti non sono molto chiare rispetto all'azione che si andrà ad effettuare;
- **Sicurezza:** I criteri di scelta della password risultano essere poco restrittivi, ciò può comportare una facile manomissione del sistema da attacchi esterni;

## Pseudo-requisiti

L'applicativo è stato sviluppato in Python come linguaggio di programmazione e per la web application è stato utilizzato django con il database di sistema mantenuto sul DBMS relazionale SQLite.



## Modello dei casi d'uso

Sono stati identificati due attori principali all'interno del sistema:

- Utente ospite
- Utente registrato

Essi hanno quasi le stesse funzionalità in comune, la differenza fondamentale è che l'utente registrato può salvare i risultati delle detection.

### UC\_01: Upload file

**Attore:** Utente ospite, Utente registrato

**Entry condition:** L'utente sta visualizzando la home page.

**Flusso di eventi:**

1. L'utente seleziona il campo "select file" dove dovrà inserire il dataset sul quale vuole effettuare l'analisi;
2. Il sistema fa visualizzare la schermata di selezione del file;
3. L'utente seleziona il file, che deve essere di tipo .csv;
4. Il sistema informa l'utente di aver ricevuto correttamente il file;
5. L'utente fa l'upload del file nel sistema;
6. Il sistema informa l'utente che l'upload è avvenuto con successo;
7. L'utente può cambiare il file oppure proseguire la detection.

**Exit condition:** L'utente ha caricato con successo il file.

**Flussi alternativi/eccezioni:**

- Se al passo 5 l'utente tenta di fare l'upload di un file con il formato errato, verrà eseguito il caso d'uso 01.01.
- Se al passo 7 l'utente decide di cambiare file, verrà eseguito il caso d'uso 01.02.

### UC\_01.01: Formato file errato

**Attore:** Utente ospite, Utente registrato

**Entry condition:** L'utente ha provato ad inserire il file sul quale verrà effettuata la detection ma non era del formato corretto.

**Flusso di eventi:**

1. Il sistema mostra un messaggio di errore "unsupported file type";
2. Il sistema ripropone la schermata di selezione del file.

**Exit condition:** L'utente carica il file del formato corretto.

## UC\_01.02: Cambio file

**Attore:** Utente ospite, Utente registrato

**Entry condition:** L'utente accede alla funzionalità di cambio file.

**Flusso di eventi:**

1. L'utente richiede di cambiare file;
2. Il sistema rimuove il file precedentemente caricato e ripropone la schermata di caricamento del file.

**Exit condition:** L'utente carica il nuovo file.

## UC\_02: Data smell detection

**Attore:** Utente ospite, Utente registrato

**Entry condition:** L'utente sta visualizzando la pagina con il dataset caricato con successo.

**Flusso di eventi:**

1. L'utente procede con la fase successiva di personalizzazione;
2. Il sistema fornisce una serie di parametri e scelte che permettono di personalizzare la detection;
3. L'utente seleziona i parametri più adatti al tipo di detection che vuole effettuare e invia le informazioni al sistema;
4. Il sistema salva le informazioni e informa l'utente di aver salvato con successo le informazioni;
5. L'utente procede con la fase successiva;
6. Il sistema mostra i risultati della detection sul dataset, divisi per colonne differenti, accessibili nuovamente dall'utente nella sezione "saved results".

**Exit condition:** L'utente ha completato la detection sul dataset e ha ricevuto i risultati.

**Flussi alternativi/eccezioni:**

- Se al passo 3 l'utente decide di non selezionare alcun parametro, verrà eseguito il caso d'uso 02.01.
- Se al passo 6 l'utente registrato decide di eliminare il risultato, verrà eseguito il caso d'uso 02.02.

#### UC\_02.01: Parametri non selezionati

**Attore:** Utente ospite, Utente registrato

**Entry condition:** L'utente non seleziona alcun data smell da individuare oppure nessuna colonna del dataset sulla quale effettuare la detection.

**Flusso di eventi:**

1. Il sistema fornisce un messaggio di errore "Select smells and columns";
2. Il sistema ripropone la schermata di personalizzazione della detection.

**Exit condition:** L'utente seleziona almeno uno smell e almeno una colonna.

#### UC\_02.02: Cancellazione risultati

**Attore:** Utente registrato

**Entry condition:** L'utente registrato decide di eliminare i risultati.

**Flusso di eventi:**

1. Il sistema chiede all'utente la conferma per eliminare i risultati;
2. L'utente dà il suo consenso;
3. Il sistema elimina i risultati, fornendo a schermo un messaggio di conferma.

**Exit condition:** I risultati vengono eliminati dall'account dell'utente registrato.

#### UC\_03: Visualizzazione risultati salvati

**Attore:** Utente registrato

**Entry condition:** L'utente visualizza l'home page.

**Flusso di eventi:**

1. L'utente seleziona il campo "saved results" per accedere ai risultati;
2. Il sistema fornisce la schermata di visualizzazione dei risultati.

**Exit condition:** L'utente visualizza la pagina con i risultati salvati.

**Flussi alternativi/eccezioni:**

- Se al passo 2 il sistema non riscontra alcun risultato, verrà eseguito il caso d'uso 03.01.

#### UC\_03.01: Nessun risultato salvato

**Attore:** Utente registrato

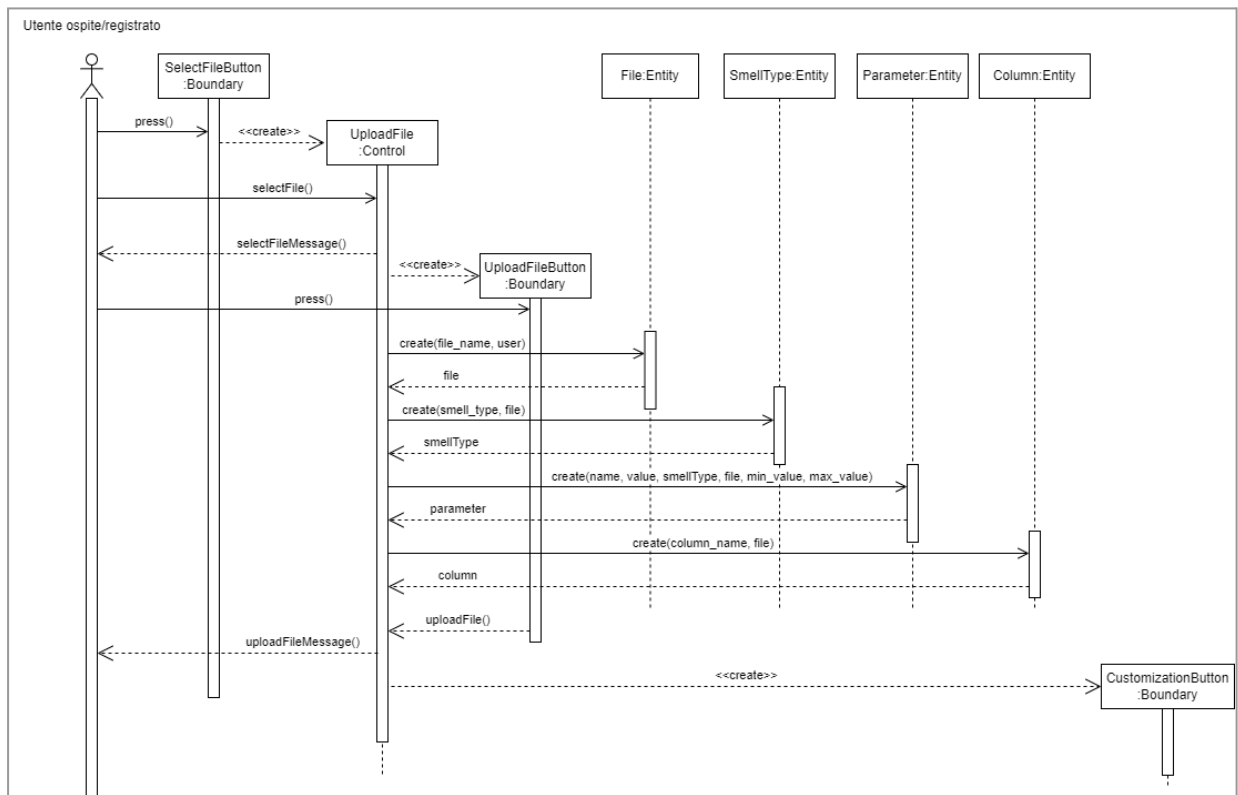
**Entry condition:** Il sistema non riscontra alcun risultato salvato.

**Flusso di eventi:** Il sistema restituisce una schermata con il messaggio "There aren't any saved detection results".

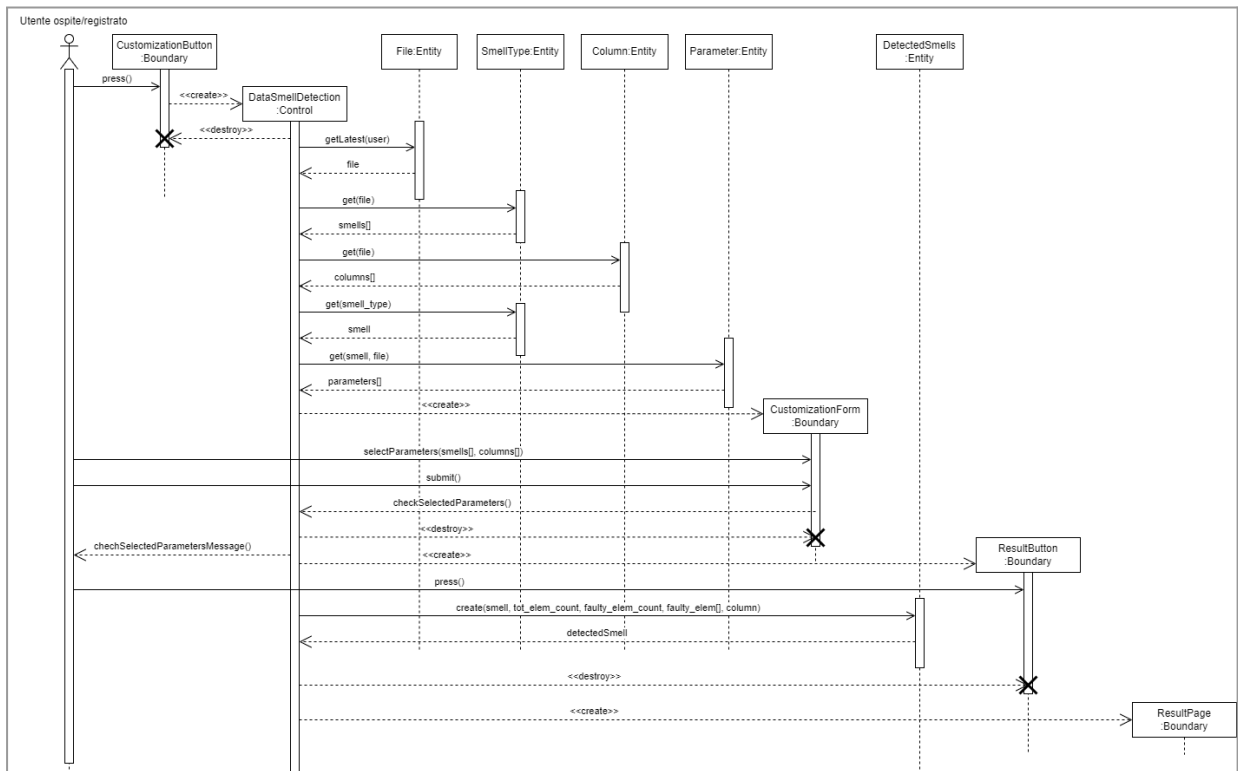
**Exit condition:** L'utente visualizza il messaggio.

# Sequence diagram

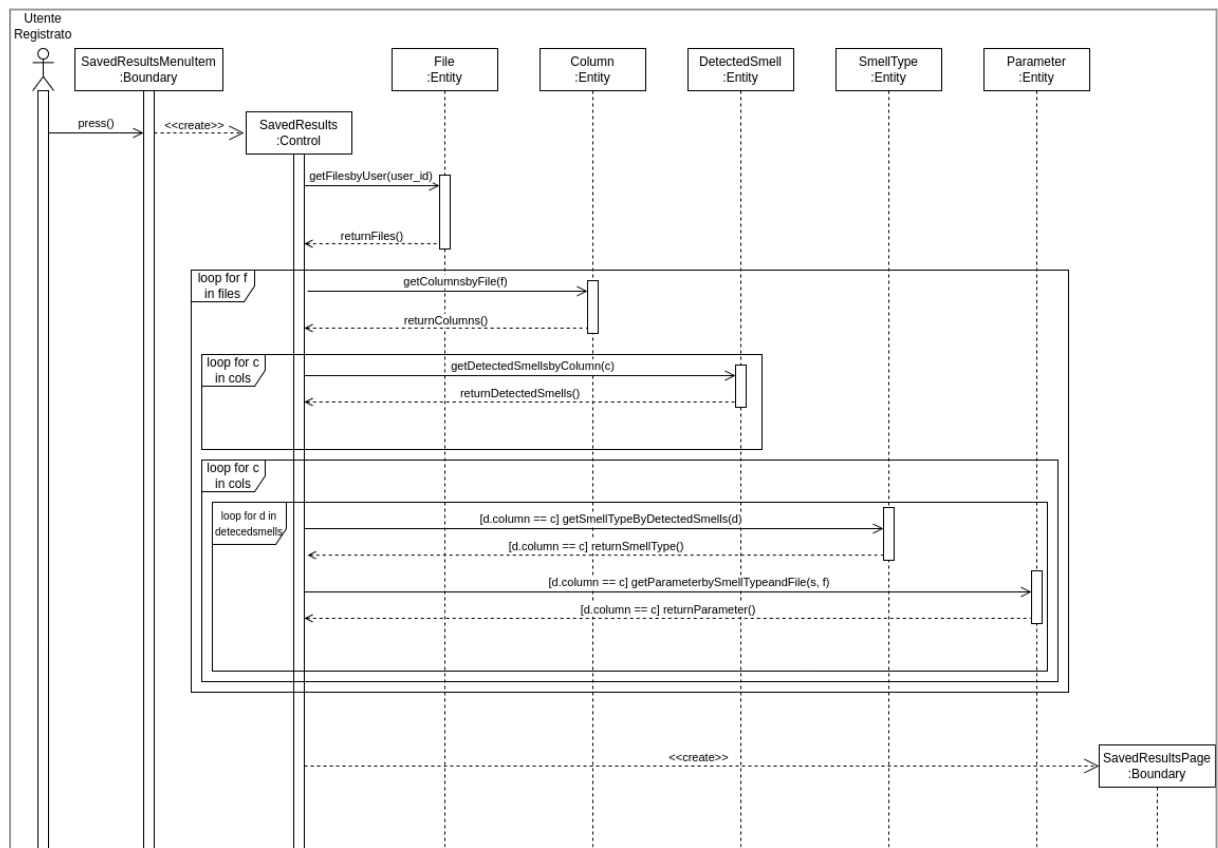
## SD\_01: Upload file



## SD\_02: Data smell detection



## SD\_03: Visualizzazione risultati salvati



# Change Requests

## Change request 1

**Change request id:** CR\_01

**Titolo:** Estensione della suite di detection di data smell.

**Descrizione:** La seguente change request prevede la realizzazione di almeno una nuova detection di data smell, in modo tale da estendere la suite di detection che mette a disposizione il tool.

**Tipologia di manutenzione:** additiva/perfettiva

**Requisiti coinvolti:** RF\_07

**Stato:** presentata/sottomessa

## Change request 2

**Change request id:** CR\_02

**Titolo:** Funzionalità di calcolo delle data quality dimensions metrics.

**Descrizione:** La seguente change request prevede la realizzazione di un meccanismo di calcolo, rispetto ai risultati ottenuti dall'analisi del tool, di un sottoinsieme di data quality dimensions (metriche che consentono di riuscire a definire la qualità dei dataset che analizza il tool DSD):

- **Completeness:** La percentuale di dati memorizzati rispetto al potenziale di "100% completo";
- **Uniqueness:** Nulla verrà registrato più di una volta in base a come viene identificato;
- **Timeliness:** Il grado di rappresentazione della realtà da parte dei dati nel momento richiesto;
- **Validity:** I dati sono validi se sono conformi alla sintassi (formato, tipo, intervallo) della loro definizione;
- **Accuracy:** Il grado in cui i dati descrivono correttamente l'oggetto o l'evento del "mondo reale" che viene descritto;
- **Consistency:** L'assenza di differenze, quando si confrontano due o più rappresentazioni di un elemento con una definizione.

**Tipologia di manutenzione:** additiva/perfettiva

**Requisiti coinvolti:** //

**Stato:** presentata/sottomessa

## Change request 3

**Change request id:** CR\_03

**Titolo:** Estensione del meccanismo di reporting.

**Descrizione:** La seguente change request prevede l'estensione dell'attuale meccanismo di reporting, aggiungendo la possibilità di visualizzare a schermo i risultati ottenuti dall'analisi di un dataset e mostrare nel tempo le differenze di risultati, tra differenti versioni di uno stesso dataset, dei relativi smell individuati e metriche/misure calcolate.

**Tipologia di manutenzione:** additiva/perfettiva

**Requisiti coinvolti:** RF\_09

**Stato:** presentata/sottomessa

# Impact Analysis

Attraverso le CR precedentemente formulate, è possibile individuare l'insieme iniziale dei moduli relativi all'implementazione del sistema che potrebbero essere oggetto di impatto, a favore delle modifiche richieste al sistema. Nella sezione [Selezione Change Request](#), è stata effettuata anche un'analisi di fattibilità sulla realizzazione delle tre CR. In questa sezione saranno invece individuati tutti i moduli coinvolti per le modifiche da effettuare. Non avendo alcuna documentazione per le fasi di raccolta e analisi requisiti, saranno effettuate delle considerazioni sull'impatto a partire dal class diagram del [Package detector](#).

## Impatto di CR\_01

Per questa change request sarà necessario aggiungere nuove istanze di detection, in particolare bisogna aggiungere dei nuovi moduli nel package ***datasmelldetection/detectors/great\_expectations/expectations***. Per una corretta visualizzazione della scelta dei nuovi smell nella [Customization Page](#) (*customization.html*), sarà necessario modificare i file ***presettings.json*** e ***smells.json*** presenti in ***web\_application/argon-dashboard-django/app***, in questo modo l'utente avrà la possibilità di scegliere i nuovi smell e di decidere al meglio i parametri da utilizzare. Le nuove detection saranno utilizzate dal detector. All'interno del modulo ***datasmelldetection/detectors/great\_expectations/detector.py*** è stata implementata la classe ***GreatExpectationsDetector*** che rappresenta la funzionalità principale del progetto. Per far sì che una detection vada a buon fine, il detector ha bisogno di istanziare le seguenti classi, alcune delle quali sono direttamente parte di questo progetto:

- ***DataContext***: Classe che imposta il contesto necessario per il corretto funzionamento della libreria di terze parti GreatExpectations;
- ***DatasetWrapper***: Classe che “wrappa” il dataset da analizzare;
- ***DataSmellRegistry***: Classe che contiene il registro delle detection di smell utilizzabili;
- ***DatasetProfiler***: Classe che associa le colonne del dataset con la relativa detection contenuta nel DataSmellRegistry;
- ***DetectionResultConverter***: Classe che si occupa di formattare i risultati della detection in un formato conveniente per la webapp;
- ***DataSmellAwareConfiguration***: Classe che memorizza la configurazione personalizzata da utilizzare per effettuare le detection.

Per l'introduzione delle nuove detection di smell, potrebbe quindi essere necessario effettuare delle modifiche alle classi descritte in precedenza.



**Starting Impact Set** = {CL\_02, CL\_03, CL\_05, CL\_06, CL\_08}.

```

graph TD
    CL_06[CL_06] --> CL_02[CL_02]
    CL_06 --> CL_05[CL_05]
    CL_06 --> CL_03[CL_03]
    CL_06 --> CL_08[CL_08]
    CL_06 --> MOD_05.2[MOD_05.2]
    CL_02 --> MOD_03.1[MOD_03.1]
    CL_05 --> MOD_04.2[MOD_04.2]
    CL_03 --> MOD_05.1[MOD_05.1]
    CL_08 --> MOD_02.2[MOD_02.2]
    MOD_05.2 --> MOD_03.1
    MOD_05.2 --> MOD_04.2
    MOD_05.2 --> MOD_05.1
    MOD_05.2 --> MOD_06.1[MOD_06.1]
    MOD_05.2 --> MOD_02.2
    MOD_03.1 --> TSC_01.2[TSC_01.2]
    MOD_04.2 --> MOD_04.1[MOD_04.1]
    MOD_04.2 --> TSC_02.1[TSC_02.1]
    MOD_04.2 --> MOD_06.1
    MOD_05.1 --> TSC_02.2[TSC_02.2]
    MOD_06.1 --> TSC_05.1[TSC_05.1]
    MOD_02.2 --> TSC_02.1
    MOD_02.2 --> TSC_05.1
  
```

[illegible][illegible]

Attraverso un approccio Distance based, si selezionano tutti gli artefatti entro tre archi di distanza a partire dalla classe principale CL\_06, ottenendo il CIS:

- **Candidate Impact Set** = {CL\_02, MOD\_03.1, TSC\_01.2, CL\_05, MOD\_04.2, MOD\_04.1, TSC\_02.2, CL\_03, MOD\_05.1, CL\_08, MOD\_02.2, MOD\_02.1, MOD\_05.2, MOD\_06.1, TSC\_05.1}.

## Impatto di CR\_02

La funzionalità descritta da questa change request non è presente nel sistema, non è possibile mappare nessun modulo esistente che possa essere modificato. In questo caso bisogna aggiungere un nuovo modulo che consenta di calcolare le data quality dimensions, i cui risultati saranno mostrati all'interno del nuovo meccanismo di reporting, che sarà anch'esso oggetto di modifica nella successiva change request. Questo modulo sarà posizionato nella cartella **app/** all'interno del package **web\_application**.

## Impatto di CR\_03

Per questa change request sarà necessario migliorare la presentazione dei risultati, in particolare bisogna modificare ciò che sarà visualizzato nella [Saved results page \(saved.html\)](#) ed eventualmente introdurre delle nuove pagine. Tutte le pagine dell'applicazione sono contenute nel package **web\_application/argon-dashboard-django/core**, alla subfolder **templates/**. Ogni volta che si aggiunge una nuova pagina al progetto, è necessario modificare i seguenti moduli:

- **web\_application/argon-dashboard-django/app/views.py**: Questo modulo si occupa di elaborare le richieste inviate alla webapp per fornire una risposta all'utente.
- **web\_application/argon-dashboard-django/app/urls.py**: Questo modulo si occupa di recuperare le richieste HTTP e di individuare un match tra l'URL richiesto e gli URL memorizzati nella webapp. Se c'è un match, viene invocata la funzione associata che è presente in *views.py*, la quale si occuperà di eseguire tutta la logica necessaria, che permetterà di fornire una risposta per mostrare la corretta visualizzazione della pagina richiesta e del relativo contesto.

Nel caso sia necessario introdurre nuove entità nel database, è necessario modificare il seguente modulo:

- **web\_application/argond-dashboard-django/app/models.py**: Questo modulo si occupa di creare le tabelle del database in SQLite tramite tecniche di ORM.

# Glossario acronimi per la tracciabilità degli artefatti

## Design

- **CL\_01** - Classe *CSVDataSetManager* del class diagram
- **CL\_02** - Classe *DatasetWrapper* del class diagram
- **CL\_03** - Classe *DataSmellAwareConfiguration* del class diagram
- **CL\_04** - Classe *DataSmellMetadata* del class diagram
- **CL\_05** - Classe *DataSmellRegistry* del class diagram
- **CL\_06** - Classe *GreatExpectationsDetector* del class diagram
- **CL\_07** - Classe *DataSmell* del class diagram
- **CL\_08** - Classe *DetectionResultConverter* del class diagram
- **CL\_09** - Classe *ExtendedDetectionResult* del class diagram
- **CL\_10** - Classe *DataSmellAwareProfiler* del class diagram

## Codice

### **PKG\_01** - Package *datasmelldetection/detectors/great\_expectations*

- **MOD\_01** - modulo *context.py*:
  - **MOD\_01.1** - class *GreatExpectationsContextBuilder*, implementazione del pattern Builder per facilitare la creazione di istanze di *DataContext*;
- **MOD\_02** - modulo *converter.py*:
  - **MOD\_02.1** - class *ExtendedDetectionResult*, implementazione della classe *CL\_09* - *ExtendedDetectionResult*;
  - **MOD\_02.2** - class *StandardResultConverter*, implementazione della classe *CL\_08* - *DetectionResultConverter* per codificare i risultati restituiti dalla libreria GreatExpectations, nel formato definito nell'implementazione della classe *ExtendedDetectionResult*;
- **MOD\_03** - modulo *dataset.py*:
  - **MOD\_03.1** - class *DatasetWrapper*, implementazione della classe *CL\_02* - *DatasetWrapper*;
  - **MOD\_03.2** - class *FileBasedDataSetManager*, implementazione della classe *CL\_01* - *CsvDataSetManager*;

- **MOD\_04 - modulo *datasmell.py*:**
  - **MOD\_04.1 - class *DataSmellMetadata***, implementazione della classe *CL\_04 - DataSmellMetadata*;
  - **MOD\_04.2 - class *DataSmellRegistry***, implementazione della classe *CL\_05 - DataSmellRegistry*;
  - **MOD\_04.3 - class *DataSmell***, implementazione della classe *CL\_07 - DataSmell*;
- **MOD\_05 - modulo *detector.py*:**
  - **MOD\_05.1 - class *DataSmellAwareConfiguration***, implementazione della classe *CL\_03 - DataSmellAwareConfiguration*
  - **MOD\_05.2 - class *GreatExpectationsDetector***, implementazione della classe *CL\_06 - GreatExpectationsDetector*
  - **MOD\_05.3 - class *DetectorBuilder***, implementazione del pattern builder per facilitare la creazione di istanze di *GreatExpectationsDetector*
- **MOD\_06 - modulo *profiler.py*:**
  - **MOD\_06.1 - class *DataSmellAwareProfiler***, implementazione della classe *CL\_10 - DataSmellAwareProfiler*

# Testing

(per TSC si intende Test Suites Collection)

- **TSC\_01 - modulo *test\_dataset.py*:**
  - **TSC\_01.1 - class *TestFileBasedDatasetManager***, test suite che implementa i test case per *MOD\_03.2 - class *FileBasedManager**
  - **TSC\_01.2 - class *TestDatasetWrapper***, test suite che implementa i test case per *MOD\_03.1 - *DatasetWrapper**
- **TSC\_02 - modulo *test\_datasmell.py*:**
  - **TSC\_02.1 - class *TestDataSmellMetadata***, test suite che implementa i test case per la classe *MOD\_04.1 - class *DataSmellMetadata**
  - **TSC\_02.2 - class *TestDataSmellRegistry***, test suite che implementa i test case per la classe *MOD\_04.2 - class *DataSmellRegistry**
  - **TSC\_02.3 - class *TestDataSmell***, test suite che implementa i test case per la classe *MOD\_04.3 - *DataSmell**
- **TSC\_03 - modulo *test\_detector.py*:**
  - **TSC\_03.1 - class *TestDetectorBuilder***, test suite che implementa i test case per la classe *MOD\_05.3 - *DetectorBuilder**
- **TSC\_04 - modulo *test\_expectations.py*:**
  - **TSC\_04.1 - class *TestExpectations***, test suite che implementa i test per *PKG\_02 - *datasmelldetection/detectors/great\_expectations/expectations**
- **TSC\_05 - modulo *test\_profiler.py*:**
  - **TSC\_05.1 - class *TestDataSmellAwareProfile***, test suite che implementa i test per *MOD\_06.1 - class *DataSmellAwareProfiler**
- **TSC\_06 - modulo *test\_registration.py*:**
  - **TSC\_06.1 - class *TestExpectationRegistration***, test suite che implementa i test per la registrazione delle *Expectations* del *PKG\_02* in una istanza della classe *DataSmellRegistry*

# Selezione Change Request

Le change request sono state prioritizzate in base all'effort che è stato determinato nel realizzare le determinate modifiche e all'importanza che hanno.

Ecco di seguito la lista che mostra l'ordinamento delle CR:

1. **CR\_02 - Funzionalità di calcolo delle data quality dimensions metrics:**

Questa CR è quella che tocca meno moduli all'interno del sistema, poiché si tratta semplicemente di inserire delle funzionalità di calcolo matematico sulle colonne dei dataset. L'effort richiesto dunque è molto basso e, inoltre, senza la CR\_02 non è possibile avere la realizzazione completa della CR\_03, poiché tra le informazioni nuove che saranno aggiunte all'interno del meccanismo di reporting, ci saranno anche queste metriche. Le metriche che si intendono calcolare sono:

- **Completeness**
- **Uniqueness**
- **Validity**

2. **CR\_03 - Estensione del meccanismo di reporting:** Questa CR consente di inserire all'interno del sistema un meccanismo di reporting aggiornato, che permette di avere informazioni dettagliate sul perché quel dataset soffre di quei data smell e avere una tracciabilità nel tempo di come le modifiche al dataset possano migliorarlo o peggiorarlo, in base alle metriche calcolate in precedenza. Risulta essere la funzionalità principale che si vuole introdurre e dunque non avrebbe senso posticiparla.

3. **CR\_01 - Estensione della suite di detection di data smell:** Questa CR consente di estendere la suite di detection di data smell, andando ad inserire due nuove detection:

- **Spacing**
- **Special Character**

Nel dettaglio, questa modifica è stata lasciata come ultima poiché, richiede una approfondita comprensione del funzionamento della libreria del detector che può richiedere un maggiore sforzo, il quale potrebbe rallentare di molto il processo.

# Forward Engineering

Prima di vedere nel dettaglio il processo specifico per la realizzazione delle CR, viene fornita una breve panoramica di alcune correzioni effettuate al sistema per poter implementare le CR individuate:

- **Rimozione Utente guest:** L'utente guest una volta effettuate le analisi e visualizzati i risultati, non aveva alcun modo per interagire e cancellare i risultati ottenuti e, inoltre, per il tipo di sistema che si aveva a disposizione, l'accesso come guest non era coerente con quello che è lo scopo del tool. L'utente guest è stato rimosso, dando accesso al tool solo ad utenti registrati, che devono effettuare il login per utilizzarlo.
- **Correzione upload del file:** L'applicazione permetteva di effettuare l'upload di file .csv vuoti o che presentavano solamente l'intestazione delle colonne. Questa problematica è stata risolta.
- **Cancellazione fisica del file del dataset analizzato:** Dopo che l'utente registrato effettua un'analisi, la cancellazione dei risultati ottenuti rimuove esclusivamente le informazioni memorizzate sul database e il "riferimento logico" del file del dataset analizzato, senza però andarlo a rimuovere fisicamente dal sistema. Questa problematica è stata risolta.
- **Aggiunta gruppi per il nuovo reporting system:** Sono stati aggiunti i gruppi. Il sistema consente di inserire, al momento del caricamento del file, il nome del gruppo e di conseguenza aggiungere il file a quel gruppo specifico per effettuare analisi più dettagliate di dataset specifici. I motivi per cui si è deciso di aggiungere il concetto di gruppo sono principalmente due:
  - Confronto tra più dataset di ambito simile, per stabilire rispetto agli smell individuati e le Data Quality Dimensions calcolate, quale sia il dataset più adeguato da utilizzare per un task specifico.
  - Confronto temporale tra più dataset per verificare se una modifica effettuata al dataset ha migliorato o peggiorato il medesimo, rispetto agli smell individuati e le Data Quality Dimensions calcolate.

# Analisi dei requisiti

In questa sezione, evidenzieremo i casi d'uso individuati precedentemente nel sistema, i quali hanno avuto un cambiamento dovuto dalla realizzazione delle CR.

## Modello dei casi d'uso

Come detto in precedenza, ora l'attore del sistema è unico ed è l'Utente registrato.

### UC\_CR\_01: Upload File

**Attore:** Utente registrato

**Entry condition:** L'utente sta visualizzando la home page.

**Flusso di eventi:**

1. L'utente seleziona il campo "select file" dove dovrà inserire il dataset sul quale vuole effettuare l'analisi;
2. Il sistema fa visualizzare la schermata di selezione del file;
3. L'utente seleziona il file, che deve essere di tipo .csv;
4. Il sistema informa l'utente di aver ricevuto correttamente il file e mostra la parte del form dove deve essere inserito il nome del gruppo di cui il dataset farà parte;
5. L'utente inserisce il nome del gruppo di cui il dataset farà parte e fa l'upload del file nel sistema;
6. Il sistema informa l'utente che l'upload è avvenuto con successo;
7. L'utente può cambiare il file oppure proseguire la detection.

**Exit condition:** L'utente ha caricato con successo il file.

**Flussi alternativi/eccezioni:**

- Se al passo 5 l'utente tenta di fare l'upload di un file con il formato errato, verrà eseguito il caso d'uso 01.01.
- Se al passo 5 l'utente tenta di fare l'upload di un file con il campo del form vuoto, verrà eseguito il caso d'uso 01.03.
- Se al passo 7 l'utente decide di cambiare file, verrà eseguito il caso d'uso 01.02.



#### UC\_CR\_01.01: Formato file errato

**Attore:** Utente registrato

**Entry condition:** L'utente ha provato ad inserire il file sul quale verrà effettuata la detection ma non era del formato corretto.

**Flusso di eventi:**

1. Il sistema mostra un messaggio di errore "unsupported file type";
2. Il sistema ripropone la schermata di selezione del file.

**Exit condition:** L'utente carica il file del formato corretto.

#### UC\_CR\_01.02: Cambio file

**Attore:** Utente registrato

**Entry condition:** L'utente accede alla funzionalità di cambio file.

**Flusso di eventi:**

1. L'utente richiede di cambiare file;
2. Il sistema rimuove il file precedentemente caricato e ripropone la schermata di caricamento del file.

**Exit condition:** L'utente carica il nuovo file.

#### UC\_CR\_01.03: Campo form gruppo vuoto

**Attore:** Utente registrato

**Entry condition:** L'utente ha provato ad inserire il file sul quale verrà effettuata la detection ma non ha inserito alcun testo all'interno del form.

**Flusso di eventi:**

1. Il sistema mostra un messaggio "Compila questo campo";
2. Il sistema ripropone la schermata di caricamento del file.

**Exit condition:** L'utente carica il nuovo file.

### UC\_CR\_03: Visualizzazione risultati salvati

**Attore:** Utente registrato

**Entry condition:** L'utente visualizza l'home page.

**Flusso di eventi:**

1. L'utente seleziona il campo "saved results" per accedere ai risultati delle detection salvati;
2. Il sistema fornisce la lista dei gruppi di dataset creati;
3. L'utente accede ad un gruppo specifico;
4. Il sistema mostra le metriche per il gruppo di dataset e i dataset analizzati;
5. L'utente seleziona il dataset specifico contenuto nel gruppo;
6. Il sistema mostra i risultati specifici del dataset analizzato e le metriche calcolate.

**Exit condition:** L'utente visualizza la pagina con i risultati salvati di un gruppo specifico.

**Flussi alternativi/eccezioni:**

- Se al passo 2 il sistema non riscontra alcun risultato, verrà eseguito il caso d'uso 03.01.

#### UC\_CR\_03.01: Nessun risultato salvato

**Attore:** Utente registrato

**Entry condition:** Il sistema non riscontra alcun risultato salvato.

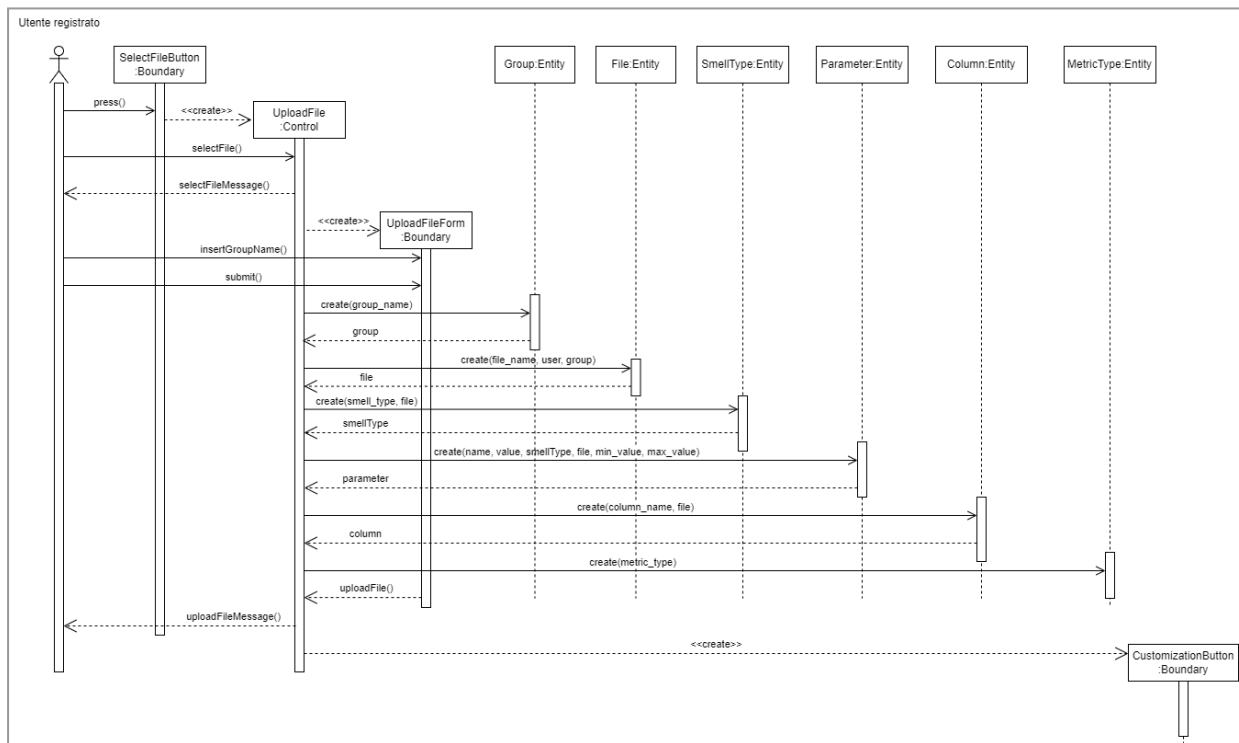
**Flusso di eventi:**

1. Il sistema restituisce una schermata con il messaggio "There aren't any saved detection results".

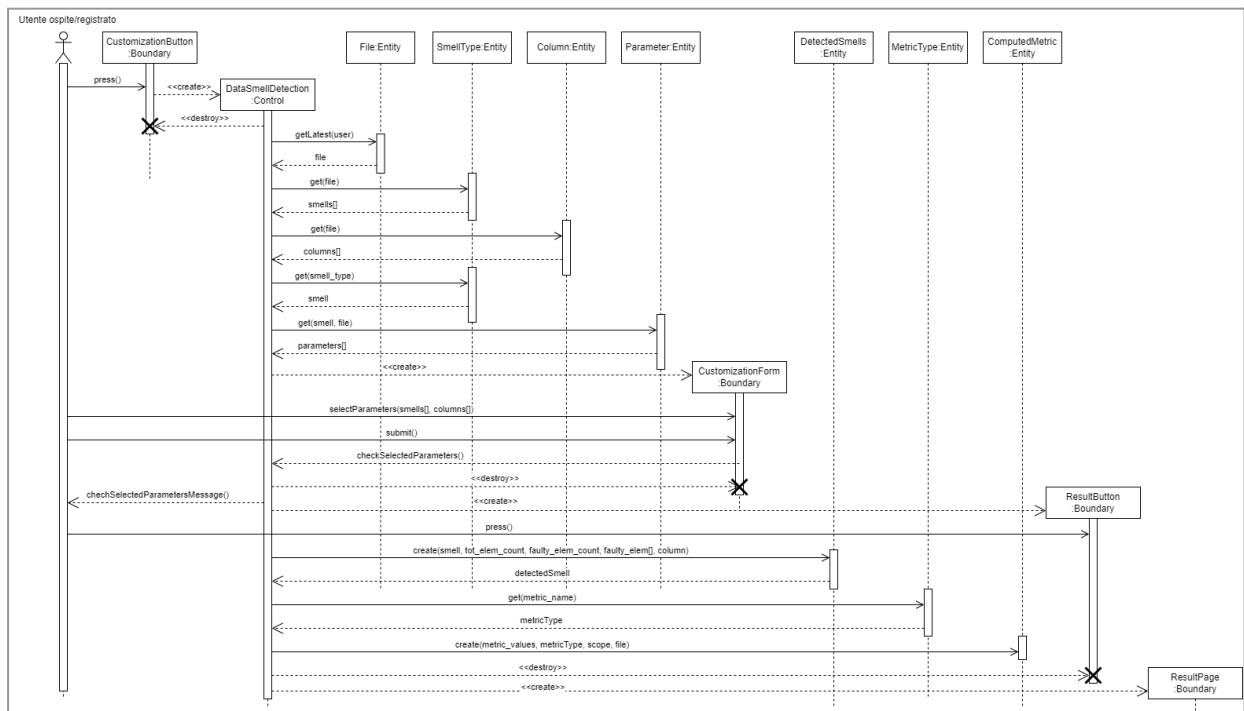
**Exit condition:** L'utente visualizza il messaggio.

# Sequence diagram

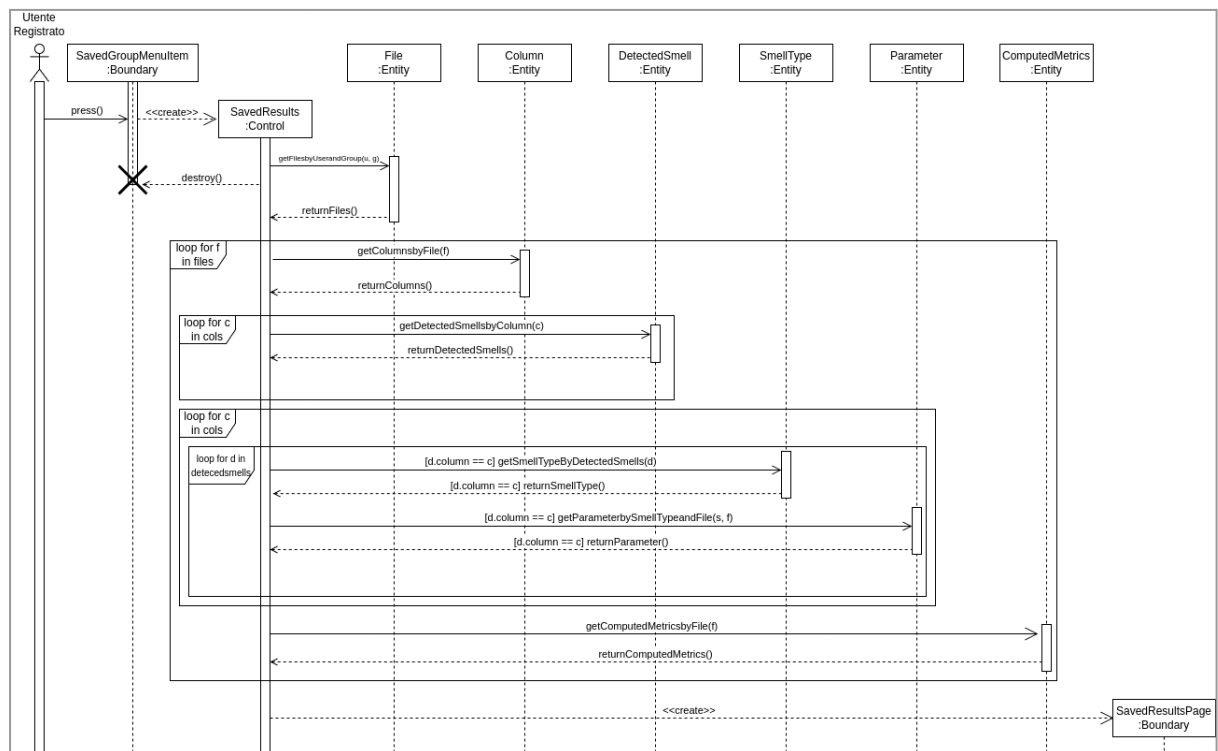
## SD\_CR\_01: Upload file



## SD\_CR\_02: Data Smell Detection



## SD\_CR\_03: Visualizzazione risultati salvati



# System & Object Design

Ci sono state poche modifiche effettuate a livello di classi presenti nel sistema. Infatti, la struttura è rimasta pressoché identica, le uniche aggiunte sono state effettuate nel package expectations, dove sono stati aggiunti due moduli per i due detector nuovi dei data smell che sono stati considerati, cioè Spacing e Special Character Smell.

## Package expectations

Per le nuove classi, implementate per gli smell specifici, è stato utilizzato un approccio regex-based, cioè sono state utilizzate delle regex, che hanno consentito di effettuare un controllo sintattico sulle colonne del dataset e capire se rispettavano o meno la regex. Di seguito mostreremo quali sono le classi aggiunte nel package già esistente:

- **ExpectColumnValuesToNotContainSpacingSmell**: Classe che consente di implementare il detector per lo Spacing Smell. Oltre al parametro *mostly*, è stato inserito un parametro *regex*, dove è stata inserita la regex che consente di individuare dei pattern specifici nelle colonne per lo smell che deve essere rilevato. La regex utilizzata è la seguente: `'^\s|\s\s+|\s$'`.
- **ExpectColumnValuesToNotContainSpecialCharacterSmell**: Classe che consente di implementare il detector per lo Special Character Smell. Anche qui è stato aggiunto il parametro *regex*. La regex utilizzata è la seguente: `'[^a-zA-Z0-9\s]'`.

## Package web\_application

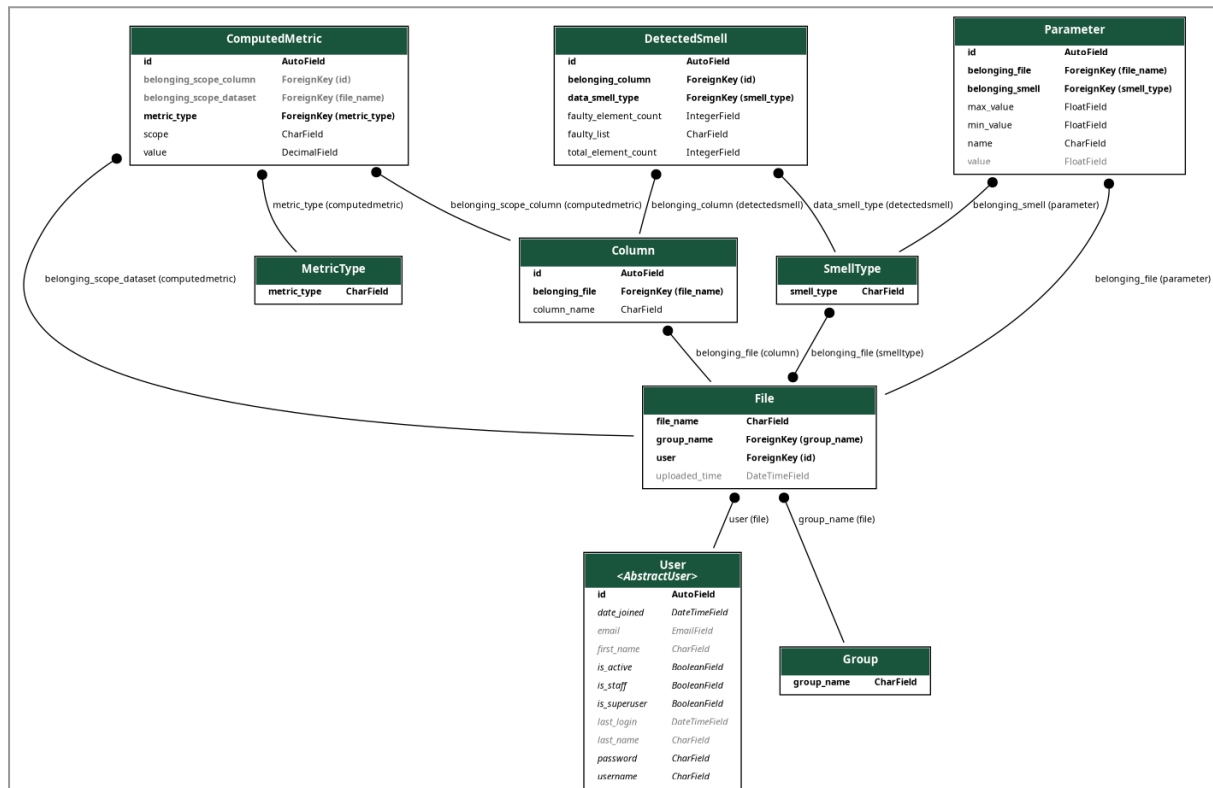
In questo package (nel dettaglio nella subfolder *app*), è stato aggiunto un nuovo modulo denominato *manage\_metrics.py*, che consente di gestire le metriche. In esso, sono state inserite le seguenti funzioni:

- **compute\_metric**: Funzione che consente di calcolare le metriche di completeness, uniqueness e validity di un dataset. Il seguente metodo è un metodo parametrico, quindi, a seconda del tipo di metrica che deve essere calcolato, ci sono parametri specifici che vengono settati, per consentire il corretto calcolo delle metriche.
- **save\_metric**: Funzione che consente di salvare le metriche nel database, distinguendo se il tipo di metrica è calcolata a livello globale oppure a livello di colonna del dataset.
- **retrieve\_metric**: Funzione che consente di prelevare le metriche calcolate precedentemente dal database, distinguendo se il tipo di metrica è calcolata a livello globale oppure a livello di colonna del dataset.

In questo modo si evita di dover (ri)calcolare le metriche ogni volta che vengono visualizzati i risultati, ma di poterle invece rapidamente recuperare dal database. Le metriche fanno parte dei risultati da mostrare all'utente rispetto al dataset analizzato. Una volta ottenute le informazioni sulla detection degli smell, le metriche saranno calcolate per una sola volta rispetto al dataset analizzato e successivamente memorizzate nel database.

## Gestione dei dati persistenti

Nella figura seguente viene mostrato il database del sistema con le modifiche effettuate su di esso per rispettare i cambiamenti delle CR implementate.



I nuovi oggetti inseriti all'interno del db sono:

- **Group**: In esso è possibile salvare il nome del gruppo che viene associato ai file, per garantire la tracciabilità nel tempo di più dataset e dunque visualizzare i dati di più dataset in un singolo gruppo.
- **MetricType**: In esso vengono salvati i tipi di metriche che sono stati inseriti nel sistema (Completeness, Uniqueness, Validity). Ciò è stato fatto per garantire un salvataggio delle informazioni in maniera persistente e per garantire una modularità, che consenta a altri developer di poter inserire nuove metriche in futuro, in maniera più chiara e semplice.
- **ComputedMetric**: In esso sono salvate tutte le informazioni che riguardano il tipo di metrica calcolato su uno specifico dataset o colonna di dataset, con il suo valore.

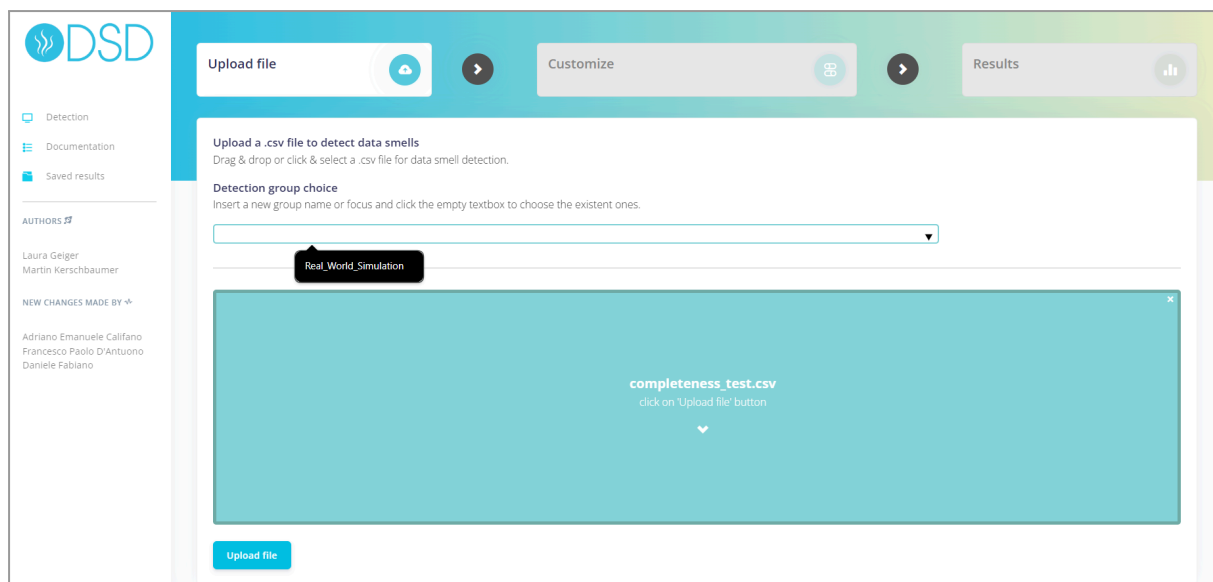
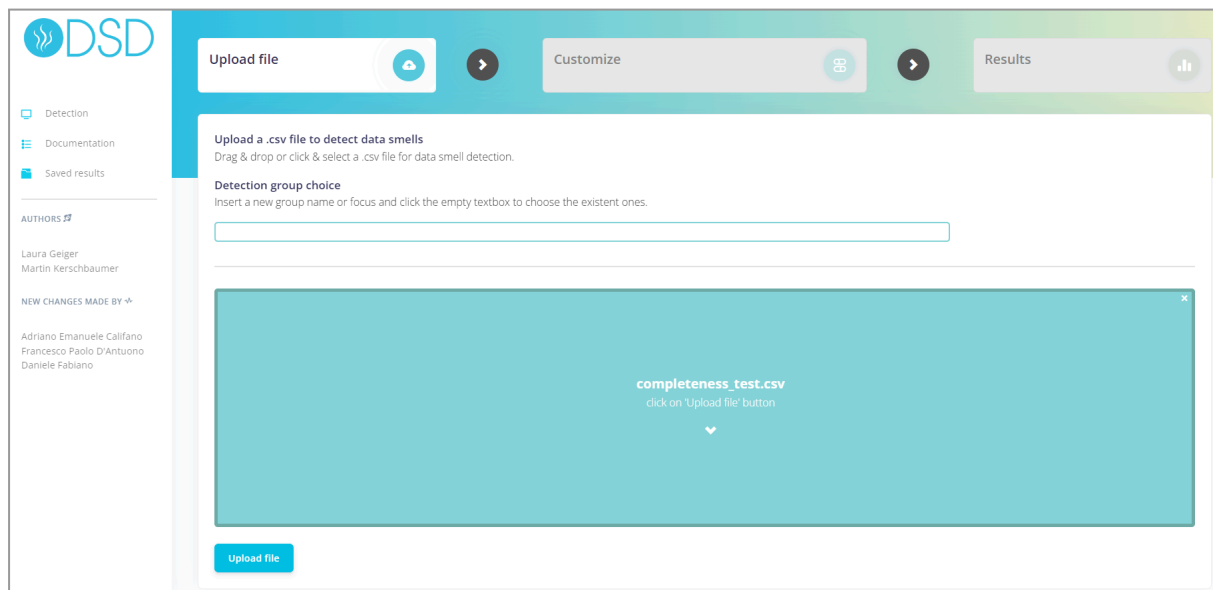
## Testing

Il testing delle nuove funzionalità del sistema è descritto nei documenti che sono specificatamente dedicati al testing.

# Descrizione delle nuove pagine

## Home page

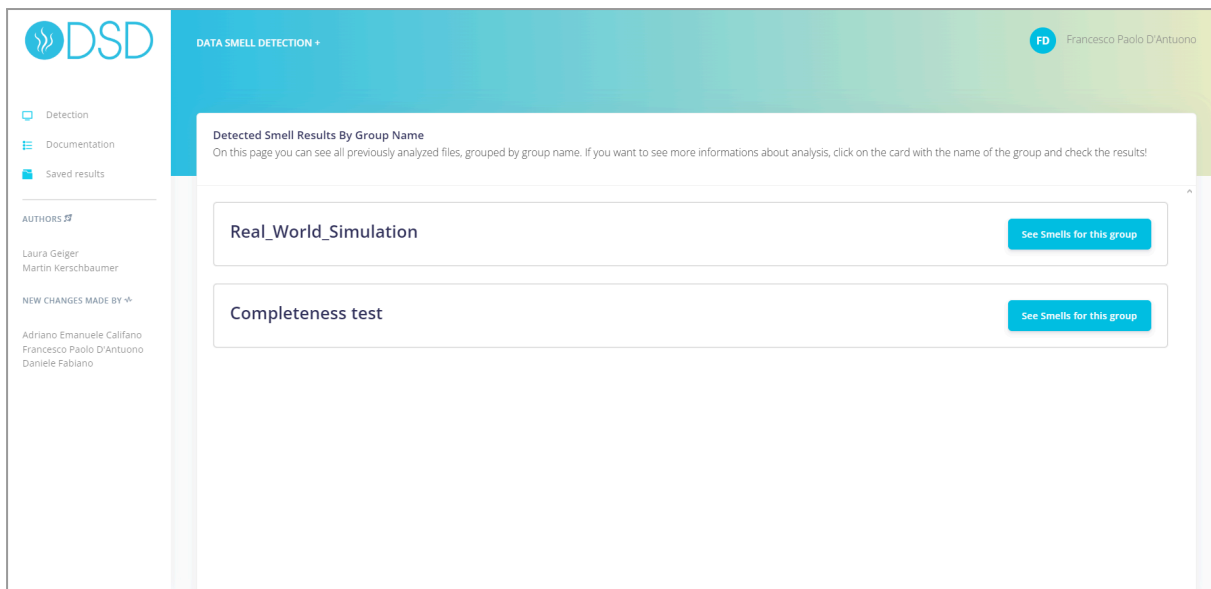
La nuova Home Page permette l'aggiunta del gruppo di risultati, dopo aver scelto il file da voler analizzare. Inoltre cliccando sulla casella di testo è possibile scegliere un gruppo già esistente dalla lista che compare.





# Saved results page

La nuova Saved results page permette di avere una visualizzazione iniziale dei vari gruppi, che sono stati creati nella fase iniziale della detection, nei quali sono presenti tutti i risultati di uno stesso gruppo. Se si vogliono avere informazioni dettagliate di un gruppo, bisogna cliccare sul pulsante alla destra della lista. La pagina del gruppo permette di visualizzare le metriche di completeness, uniqueness e validity globali. Se si vogliono avere più dettagli per la detection del singolo dataset che fa parte di quel gruppo, bisogna cliccare sul pulsante con su scritto il nome del dataset del quale si vogliono visualizzare i risultati.



In questa specifica sezione, che si potrà visualizzare solo dopo aver cliccato sul pulsante, si può avere informazione sugli smell che sono stati detectati su ogni colonna e le metriche specifiche per quel dataset.

