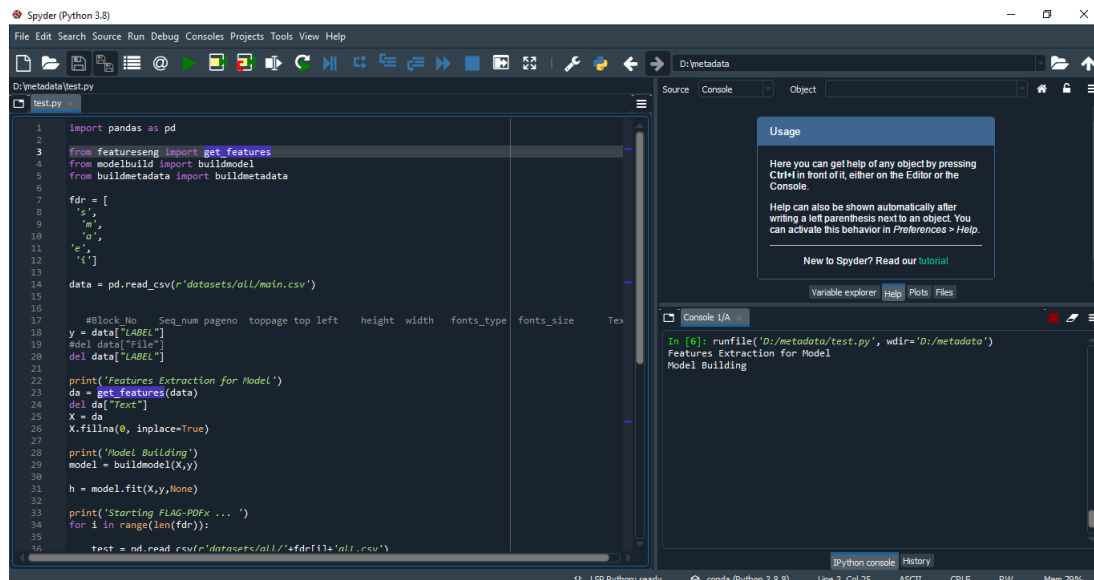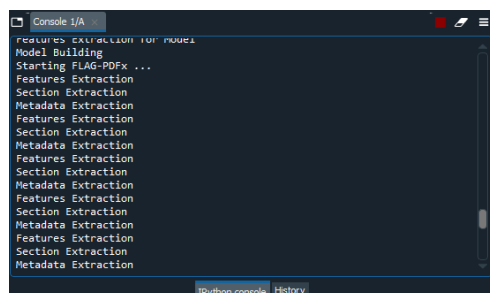Help File to Execute the code and verify extracted Metadata

1. The folders "\metadata\datasets" has the code and extracted PDF content in form of CSV files
2. Please install python 3.7 version
3. In \metadata folder is test.py file. This is the main file which run the complete code



4. When test.py is executed then console file shows the process stages. The Model Building stage is based on system specifications.
5. When complete process is completed then the system generates the output file in form of metadata CSV files.



6. The main.csv is used to build the model
7. The metadata\datasets\all has csv with file names like ialL (i, e, m, s, a represents different publishers). These files contain the extracted pdf content with geometric and fonts info.
8. The output of system is in metadata\datasets\ folder (with initial class label and section postfix).
   a. TI -    Title
   b. AU-    Authors
   c. AFF-   Affiliations
   d. EM-    Emails
   e. HL-    Header/ footers
   f. H1-    Heading level 1
   g. H2-    Heading level 2
   h. H3-    Heading level 3

       i.    TA-     Table

       j.    FI-     Figures

       k.    ACK-    Acknowledgements

       l.    RE-     References

9. Dataset Folders ds\Annotated\ has subfolders publisher wise.

10. Each csv file has the metadata of same file pdf

11. To verify output

File  Home  Insert  Page Layout  Formulas  Data  Review  View  Help  Tell me what you want to do  Share

Get External Data | New Query | Recent Sources | Refresh All | Connections | Properties | Edit Links | Sort | Filter | Clear | Reapply | Advanced | Text to Columns | Flash Fill | Remove Duplicates | Data Validation | Consolidate | Relationships | What-If Analysis | Forecast Sheet | Group | Ungroup | Subtotal

Get & Transform | Connections | Sort & Filter | Data Tools | Forecast | Outline

A1 : filename

| | A |
| --- | --- |
| 1 | filename |
| 2 | m10 |
| 3 | m10 |
| 4 | m10 |
| 5 | m10 |
| 6 | m10 |
| 7 | m10 |
| 8 | m10 |
| 9 | m10 |
| 10 | m11 |
| 11 | m11 |
| 12 | m11 |
| 13 | m11 |
| 14 | m11 |
| 15 | m11 |
| 16 | m11 |
| 17 | m12 |

Sort A to Z
Sort Z to A
Sort by Color
Sheet View
Clear Filter From "filename"
Filter by Color
Text Filters
m1
No matches
OK  Cancel

Tensile Strength (MPa) Tensile Strain Capacity (%) Stress Performance Index

Ready  35 of 2

---

File  Home  Insert  Page Layout  Formulas  Data  Review  View  Help  Tell me what you want to do

Paste | Calibri | 11 | B I U | Font | Clipboard

File  Home  Insert  Page Layout  Formulas  Data  Review  View  Help  Tell me what you want to do

Get External Data | New Query | Show Queries | From Table | Recent Sources | Refresh All | Connections | Properties | Edit Links | Sort | Filter | Clear | Reapply | Advanced | Text to Columns | Flash Fill | Remove Duplicates | Data Validation | Consolidate | Relationships

Get & Transform | Connections | Sort & Filter | Data Tools

F17 :

| | A | B | C |
| --- | --- | --- | --- |
| 1 | heading | level | |
| 2 | 1. Introdution | h1 | |
| 3 | 2. Materials and Methods | h1 | |
| 4 | 2.1. Chemicals | h2 | |
| 5 | 2.2. Samples and Reagents Preparation | h2 | |
| 6 | 2.3. CE-LIF | h2 | |
| 7 | 3. Results | h1 | |
| 8 | 3.1. Spikes Are Present in Samples of Incubated | h2 | |
| 9 | 3.2. Analytical Instrument Can Generate Artefact | h2 | |
| 10 | 3.3. Changes in Injection Volume Reveal That Sp | h2 | |
| 11 | 3.4. The Nature of Running Buffer Affects Spikes | h2 | |
| 12 | 3.5. The Link between Spikes and Aggregates Is S | h2 | |
| 13 | 3.6. The Use of Thioflavine T in Running Buffer N | h2 | |
| 14 | 3.7. The CE-LIF Separation of Aβ Aggregates Exhi | h2 | |
| 15 | 3.8. This CE-LIF Method Allows Separation of Aβ | h2 | |
| 16 | 4. Discussion | h1 | |
| 17 | Supplementary Materials: | h1 | |
| 18 | Acknowledgments | h1 | |
| 19 | | | |
| 20 | | | |
| 21 | | | |

papermetadata | authors | keywords | heads

A2 : m1

| | A | B | C | D |
| --- | --- | --- | --- | --- |
| 1 | filename | top | Text | |
| 2 | m1 | | 1551 | 2.1. Chemicals |
| 3 | m1 | | 1621 | 2.2. Samples and Reagents Preparation |
| 4 | m1 | | 5623 | 3.4. The Nature of Running Buffer Affects Spikes Detection |
| 5 | m1 | | 5800 | 3.5. The Link between Spikes and Aggregates Is Supported by Ultrasound-Treatment of Samples |
| 6 | m1 | | 5897 | 3.6. The Use of Thioﬂavine T in Running Buffer Must Be Carefully Controlled and Validated |
| | | | 321 | |
| | | | 322 | |
| | | | 323 | |
| | | | 324 | |
| | | | 325 | |
| | | | 326 | |
| | | | 327 | |
| | | | 328 | |
| | | | 329 | |
| | | | 330 | |
| | | | 331 | |

mH2

Filter Mode

**Window 1 — m1 - Excel**

| | A | B |
|---|---|---|
| 1 | title | Advances and Pitfalls in the Capil... |
| 2 | journal | MDPI Separations |
| 3 | accepted on | 14/12/2017 |
| 4 | published on | 22/12/2017 |
| 5 | pages | 16 |
| 6 | doi | 10.3390/separations5010002 |
| 7 | issue | |
| 8 | vol | |
| 9 | year | 2018 |
| 10 | reference | 33 |

Sheet tabs: **papermetadata** | authors | keywor...

Cell B10: 33

**Window 2 — mHL - Excel**

Cell C2: Separations 2018 , 5 , 2; doi:10.3390/separations5010002 www.mdpi.com/journal/separations

| | A | B | C |
|---|---|---|---|
| 1 | filenam | top | Text |
| 2 | m1 | 843 | Separations 2018 , 5 , 2; doi:10.3390/separations5010002 www.mdpi.com/journal/separations |