# Leveraging a Python IDE for Data Science

**Xavier Morera**

HELPING DEVELOPERS UNDERSTAND SEARCH & BIG DATA
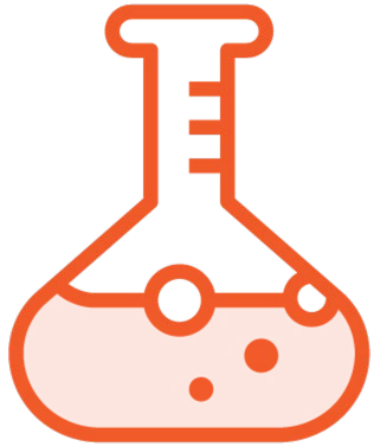
@xmorera    www.xaviermorera.com
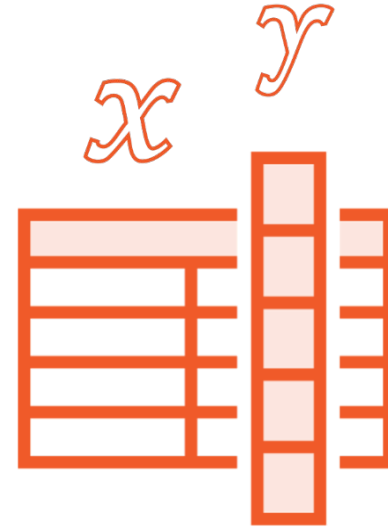
# Python IDE for Data Science

**PyCharm Pro**

**Spyder**

# Python Notebooks for Data Science

**Special kind of IDE**

**Web application**

**Execute Python code**
- Other available interpreters
- Scala, R, Shell, SQL, Spark...

# Python Notebooks for Data Science

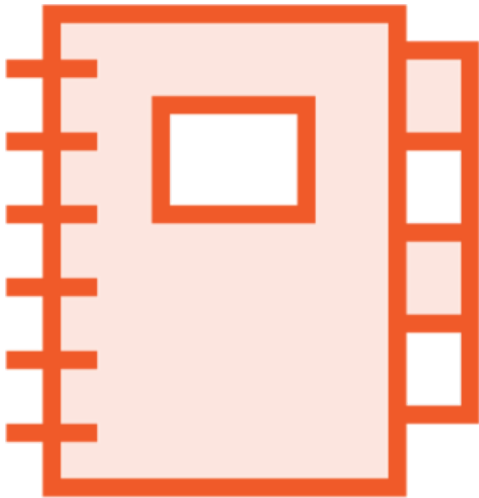**Work with data**

- Visualize

- Collaborate

**Great for iteration, prototyping , and testing**

# Python Notebooks for Data Science



**Jupyter Notebook**

**Apache Zeppelin**

# IPython Kernel-based IDE

**Cloudera Data Science
Workbench**

# Scientific Mode in PyCharm

**Support for**

- Interactive scientific computing
- And data visualization
- Requires numpy and matplotlib

# Scientific Mode in PyCharm
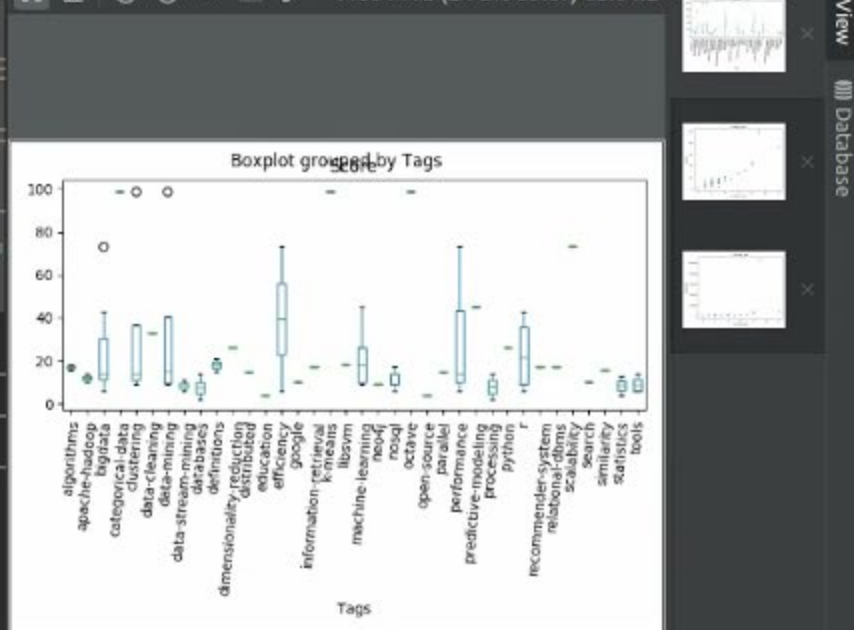
**Provides**

- IPython console

- Documentation tool window

- DataFrame/Array view

- SciView

File  Edit  View  Navigate  Code  Refactor  Run  Tools  VCS  Window  Help

ps-ide-ds ⟩ ps_ide_ds_pycharm.py                                        ps_ide_ds_pycharm ▾  ▶ 🐞 🔁 ⏱ ☰ ■  🔍

**Project** ▾  ⊕ ⊝ ✿ —    main.py ×    ps_ide_ds_pycharm.py ×      SciView:  Data  Plots            ✿ —

```
 1    # -*- coding: utf-8 -*-
 2    """
 3    Programming Python using an IDE Course at Pluralsight Data Science Demo
 4    """
 5
 6
 7    import pandas as pd
 8    def splitDataFrameList(df, target_column, separator):
 9        ''' df = dataframe to split,
10        target_column = the column containing the values to split
11        separator = the symbol used to perform the split
12        returns: a dataframe with each entry for the target column separated, with each element moved into
13        The values in the other columns are duplicated across the newly divided rows.
14        '''
15        def splitListToRows(row, row_accumulator, target_column, separator):
16            split_row = row[target_column].split(separator)
17            for s in split_row:
18                new_row = row.to_dict()
19                new_row[target_column] = s
20                row_accumulator.append(new_row)
21        new_rows = []
22        df.apply(splitListToRows, axis=1, args=(new_rows, target_column, separator))
23        new_df = pd.DataFrame(new_rows)
24        return new_df
25
26
27    # Read data from file 'posts-100-header.csv'
28    data = pd.read_csv("data/posts-100-header.csv")
29    # Check the top 3 rows
30    print(data.head(3))
31
```

Project tree:
- ps-ide-ds ~/ps-ide-ds
  - ▶ data
  - models
  - notebooks
  - main.py
  - ps_ide_ds_pycharm.py
  - README.md
  - requirements.txt
  - ▶ External Libraries
  - Scratches and Consoles

SciView panel — ‹480 PNG (24-bit color) 63.9 kB

Boxplot grouped by Tags

**ps_ide_ds_pycharm** ×                                                        ✿ —

```
import sys; print('Python %s on %s' % (sys.version, sys.platform))
sys.path.extend(['/home/xavier/ps-ide-ds'])
Python 3.6.7 (default, Oct 22 2018, 11:32:17)
In[2]: runfile('/home/xavier/ps-ide-ds/ps_ide_ds_pycharm.py', wdir='/home/xavier/ps-ide-ds')
   Id  PostTypeId  ... FavoriteCount           ClosedDate
0   5           1  ...           1.0  2014-05-14T14:40:25.950
1   7           1  ...           1.0  2014-05-14T08:40:54.950
2   9           2  ...           NaN                      NaN

[3 rows x 12 columns]

In[3]:
```

▶ 🔳 Special Variables
▶ 🔳 clean_data = {DataFrame}   Id PostTypeId ... ...View as DataFrame
▶ 🔳 data = {DataFrame}   Id PostTypeId ... Favori...View as DataFrame
▶ 🔳 tag_separated = {DataFrame}   AnswerCount...View as DataFrame
▶ 🔳 x = {Series} 0   1.0\n1   3.0\n4   4.0\n5   0.0\n6   ...View as Series
▶ 🔳 y = {Series} 0   448.0\n1   388.0\n4   1243.0\n5   ...View as Series

🐍 Python Console    Terminal    6: TODO                                    Event Log

1:24  LF :  UTF-8  4 spaces  Python 3.6

# Spyder

**An IDE built for Scientific Python**
- IDE for Data Science

**Designed by and for**
- Scientists, engineers, and data analysts

**Powerful features**

**IPython console and variable explorer**

**Extended via Spyder Notebook, Spyder Terminal...**

Editor - /home/xavier/ps-ide-ds/spyder/ps-ide-ds-spyder.py

temp.py ✕    ps-ide-ds-spyder.py ✕

```python
1  # -*- coding: utf-8 -*-
2  """
3  Programming Python using an IDE Course at Pluralsight Data Science Demo
4
5  """
6  #%%
7  import pandas as pd
8  def splitDataFrameList(df,target_column,separator):
9      ''' df = dataframe to split,
10     target_column = the column containing the values to split
11     separator = the symbol used to perform the split
12     returns: a dataframe with each entry for the target column separated, with each element moved into a new row.
13     The values in the other columns are duplicated across the newly divided rows.
14     '''
15     def splitListToRows(row,row_accumulator,target_column,separator):
16         split_row = row[target_column].split(separator)
17         for s in split_row:
18             new_row = row.to_dict()
19             new_row[target_column] = s
20             row_accumulator.append(new_row)
21     new_rows = []
22     df.apply(splitListToRows,axis=1,args = (new_rows,target_column,separator))
23     new_df = pd.DataFrame(new_rows)
24     return new_df
25  #%%
26  # Read data from file 'posts-100-header.csv'
27  data = pd.read_csv("posts-100-header.csv")
28  # Check the top 3 rows
29  print(data.head(3))
30  #%%
31  # Drop NAs from 2 columns
32  clean_data = data.dropna(subset=['AnswerCount', 'Score'])
33  # Split Tags into new rows
34  tag_separated = splitDataFrameList(clean_data, 'Tags', "><")
35  # Clean the Tags names removing '<' and '>'
36  tag_separated['Tags'] = tag_separated['Tags'].map(lambda x: x.lstrip('<').rstrip('>'))
37  #%%
38  # May take some time to display
39  import matplotlib.pyplot as plt
40  tag_separated.boxplot(by='Tags', column='Score', grid=False, rot=85)
41  plt.show()
42  #%%
43  import matplotlib.pyplot as plt
44  x = clean_data['AnswerCount']
45  y = clean_data['Score']
46
```
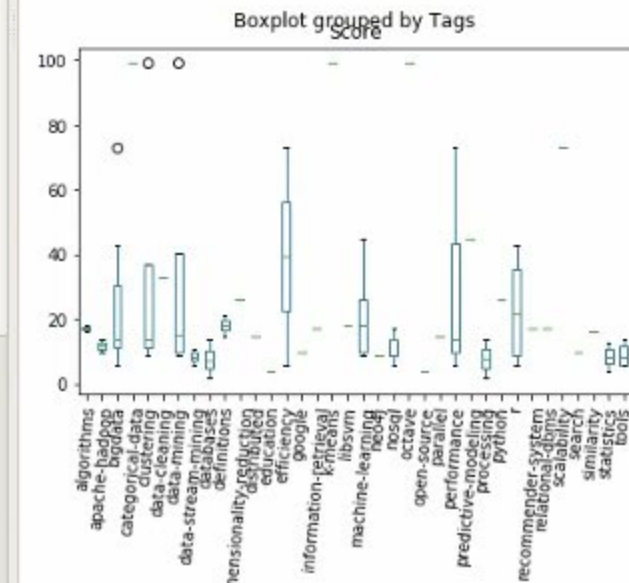
File explorer

| Name | Size | Type | Date |
|---|---|---|---|
| posts-100-header.csv | 10 KB | csv File | 5/22/ |
| ps-ide-ds-spyder.py | 1 KB | py File | 6/6/1! |

IPython console

Console 1/A ✕

```
[3 rows x 12 columns]
```



IPython console    History log

# Jupyter Notebook

**Open-source web application**
- Evolved from IPython

**Create and share documents**
- Code, equations, visualizations and text

# Jupyter Notebook
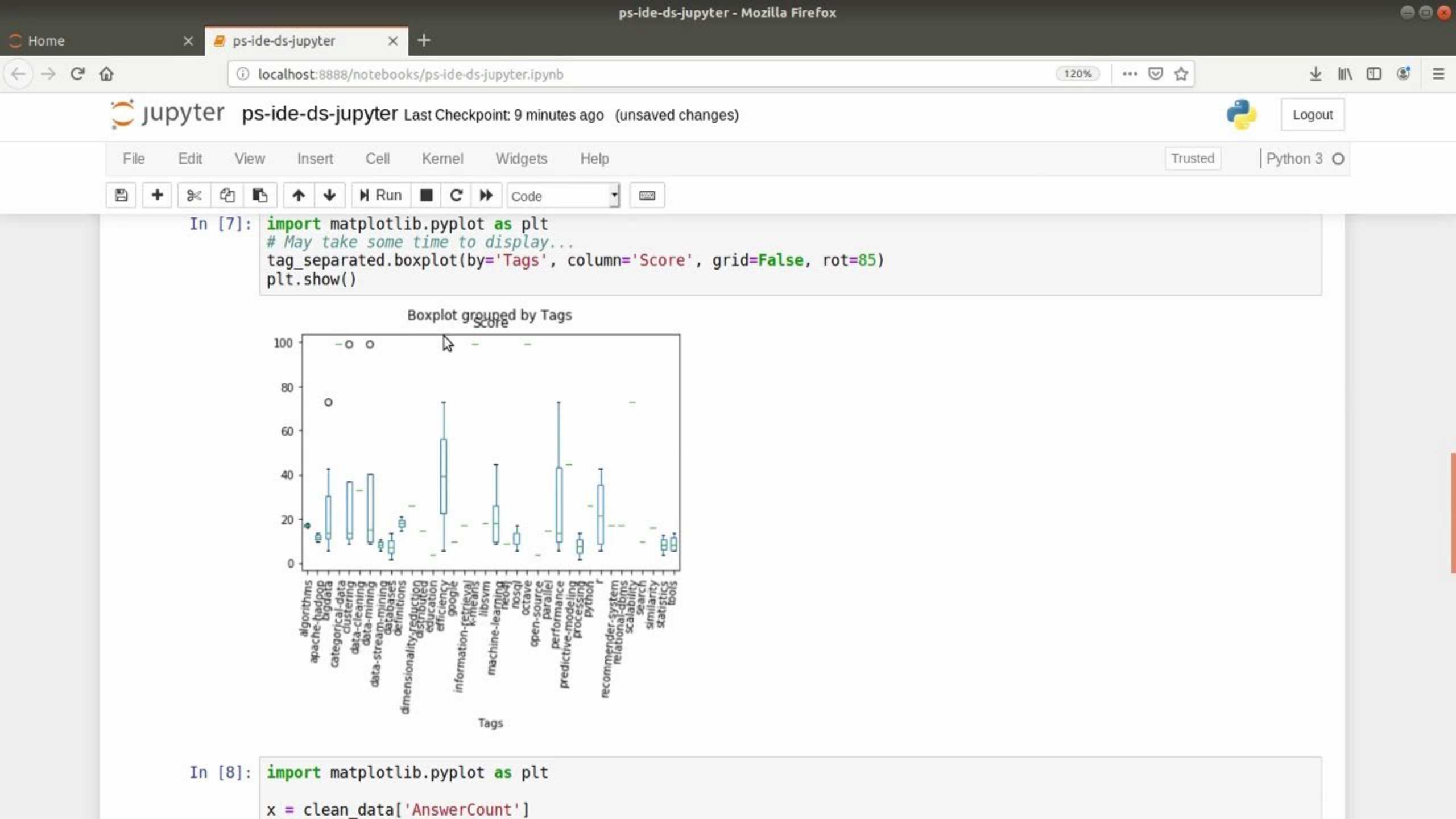
For data cleansing and transformation

Numerical simulation

Statistical modeling

Data visualization

Machine learning

jupyter ps-ide-ds-jupyter Last Checkpoint: 9 minutes ago (unsaved changes)

Logout

File   Edit   View   Insert   Cell   Kernel   Widgets   Help

Trusted | Python 3 ○

Code

```python
In [7]: import matplotlib.pyplot as plt
        # May take some time to display...
        tag_separated.boxplot(by='Tags', column='Score', grid=False, rot=85)
        plt.show()
```



Boxplot grouped by Tags
Score

```python
In [8]: import matplotlib.pyplot as plt

        x = clean_data['AnswerCount']
```

# Apache Zeppelin

**Web-based notebook**

**Data-drive interactive data analytics**

**Collaborative documents**
- Python
- SQL, Scala, shell commands...

**Single user or multi-user**

# Zeppelin

Notebook ▾    Job

Q Search

● anonymous ▾

# ps-ide-zeppelin

▷ ⋈ 📖 ✎ 🗐 ⬇    📄 ⊖ ⇌  Head ▾    Q    🗑

⌨ ⚙ 🔒  default ▾

```python
%python
import matplotlib.pyplot as plt
# May take some time to display
tag_separated.boxplot(by='Tags', column='Score', grid=False, rot=85)
plt.show()
```

FINISHED  ▷ ⋈ 📖 ⚙

# Cloudera Data Science Workbench

**IPython kernel-based IDE**

**Python development with data**
- Also Scala, and R

# Cloudera Data Science Workbench

**Platform for collaborative data science**
- At scale
- Self-service
- Leverages Git for collaboration

**Machine-learning focused**
- Quickly deploy models, with confidence
- Run experiments

```python
1   # Python Programming using an IDE Course Data Science Demos
2   import pandas as pd
3   import matplotlib as plt
4
5   # Read data from file 'posts-100-header.csv'
6   data = pd.read_csv("posts-100-header.csv")
7   # Check the top 3 rows
8   data.head(3)
9
10  def splitDataFrameList(df,target_column,separator):
11      ''' df = dataframe to split,
12      target_column = the column containing the values to split
13      separator = the symbol used to perform the split
14      returns: a dataframe with each entry for the target column separated
15      The values in the other columns are duplicated across the newly divi
16      '''
17      def splitListToRows(row,row_accumulator,target_column,separator):
18          split_row = row[target_column].split(separator)
19          for s in split_row:
20              new_row = row.to_dict()
21              new_row[target_column] = s
22              row_accumulator.append(new_row)
23      new_rows = []
24      df.apply(splitListToRows,axis=1,args = (new_rows,target_column,separ
25      new_df = pd.DataFrame(new_rows)
26      return new_df
27
28
29  # Drop NAs from 2 columns
30  clean_data = data.dropna(subset=['AnswerCount', 'Score'])
31  # Split Tags into new rows
32  tag_separated = splitDataFrameList(clean_data, 'Tags', "><")
33  # Clean the Tags names removing '<' and '>'
34  tag_separated['Tags'] = tag_separated['Tags'].map(lambda x: x.lstrip('<'
35
36  # May take some time to display...
37  tag_separated.boxplot(by='Tags', column='Score', grid=False, rot=85)
38
39  # import matplotlib.pyplot as plt
40
41  x = clean_data['AnswerCount']
42
```

Line 1, Column 1    ★    57 Lines    Python    Spaces 2

## CDSW IDE Demo ✎

🖋 Collapse    🔗 Share

By Xavier Morera — Python 3 Session — 2 vCPU / 4 GiB Memory —    `Running`
just now

## Getting Started

This is your **Python 3 session**. Your **editor** is on the left and your **input prompt** is on the bottom.

To install a package type: **!pip3 install [package_name]** at the input prompt.

To execute code from the editor, select the code and execute it with `Command-Enter` on Mac or `Ctrl-Enter` on Windows. You can also enter code at the prompt below.

Use `?command` to get help on a particular command.

>

# Takeaway

**Python IDE for Data Science**

- Notebook

**Features to enhance working with data**

- IPython console
- Execution by cell or paragraph
- Better ways to inspect data
- Embedded visualizations
- Available large amounts of compute & memory

# Takeaway

**IDE**
- PyCharm
- Spyder

**Notebook**
- Jupyter Notebook
- Apache Zeppelin

**Cloudera Data Science Workbench**