

Session 14 (4.00 PM to 6 PM)



Exploring Python's Advanced Libraries and Frameworks for Comprehensive Analysis

Dr. Vani V, Dept. of CSE,

Nitte Meenakshi Institute of Technology, Bengaluru

Outline



Setting the Context



Introduction to advanced Python libraries



High-performance data manipulation with Vaex



Interactive data exploration with D-Tale



Feature engineering with Feature Engine

Importance of Advanced Python Libraries in Enhancing Data Analysis Capabilities and Improving Workflow Efficiency

1. Efficiency and Performance:

•Handling Large Datasets:

- Advanced libraries like Vaex are designed to work with large datasets that do not fit into memory, enabling efficient out-of-core data manipulation and analysis. This capability is crucial for processing large-scale data in a timely manner.

•Optimized Computations:

- These libraries leverage optimized algorithms and data structures, significantly reducing computation time. For example, Vaex uses memory-mapped files and efficient indexing to speed up data operations.

•Parallel Processing:

- Many advanced libraries support parallel processing and distributed computing, allowing for faster data processing by utilizing multiple CPU cores or even clusters of machines.

2. Ease of Use and Flexibility:

•User-Friendly APIs:

- Libraries like D-Tale provide intuitive interfaces for data exploration and

visualization, making it easier for data professionals to interact with their data without extensive coding.

•Seamless Integration:

- These libraries integrate seamlessly with other popular Python libraries such as Pandas, NumPy, and scikit-learn, allowing for a smooth transition and interoperability between different tools in the data science ecosystem.

3. Community Support and Documentation:

•Comprehensive Documentation:

- Most advanced libraries come with extensive documentation and tutorials, providing clear guidelines on how to use their features effectively. This makes it easier for users to get started and overcome any challenges they may encounter.

•Active Community:

- An active community of users and developers contributes to continuous improvement and support. Users can seek help from community forums, GitHub issues, and other online resources, ensuring they are not alone in their learning journey.

4. Improving Workflow Efficiency:

•Streamlined Data Processing:

- Feature Engine simplifies common preprocessing tasks such as handling missing data, encoding categorical variables, and scaling features, making the data preparation phase more efficient and less error-prone.

•Interactive Analysis:

- D-Tale provides an interactive environment for data exploration, enabling quick insights and facilitating a more iterative and exploratory analysis process. This leads to faster discovery of patterns and anomalies in the data.

5. Enhancing Analytical Capabilities:

•Advanced Techniques:

- By leveraging advanced libraries, data professionals can implement sophisticated techniques and models that were previously challenging or time-consuming. This includes out-of-core algorithms, automated machine learning, and complex feature engineering.

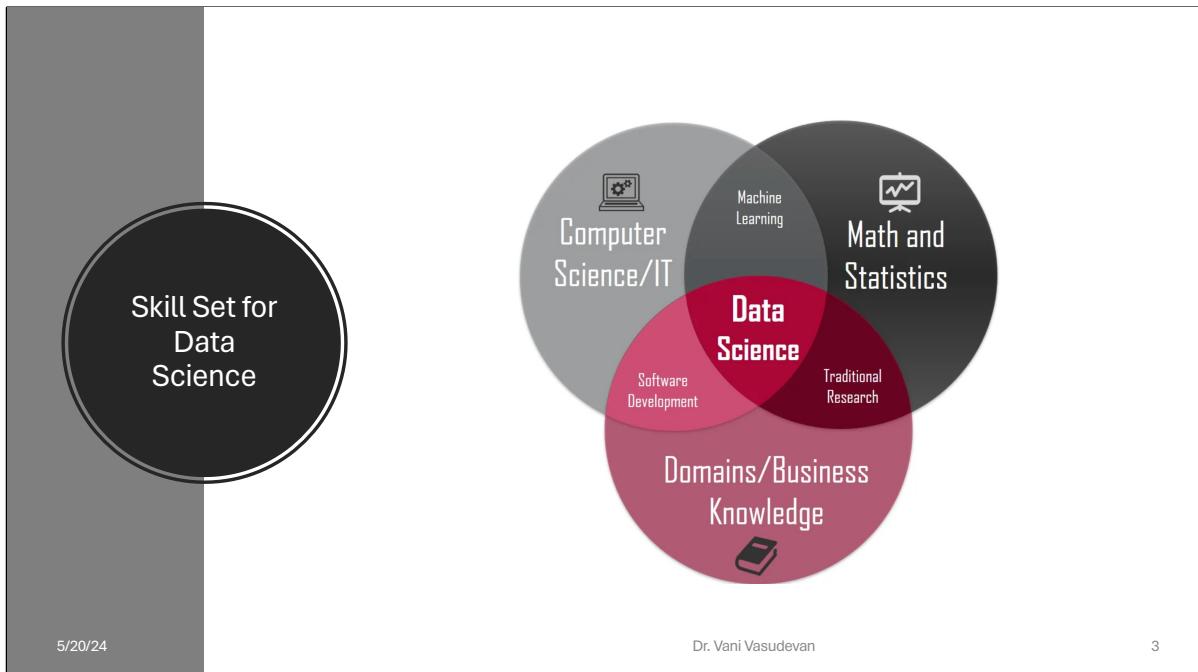
•Better Decision Making:

- The ability to analyze large datasets efficiently and accurately leads to more informed and data-driven decisions. This is particularly important in industries like healthcare, where timely and accurate insights can significantly impact business outcomes.

•Scalability:

- Advanced libraries are designed to scale with the data. As the volume and

complexity of data grow, these tools ensure that data analysis remains manageable and performant.



Data Science vs. Data Analytics

Data science is an entire field dedicated to making data more useful.

A data scientist is a professional that uses raw data to develop new ways to model data and understand the unknown.

Often, their job responsibilities incorporate **various components of computer science, predictive analytics, statistics, and machine learning**.

The collections of information that data scientists work with can be quite large, requiring expertise to organize and navigate.

Source: Google data analytics professional certificate

Data Science vs. Data Analytics

Data analytics is a subfield of the larger data science discipline.

The aim of data analytics is to create methods to capture, process, and organize data to uncover actionable insights for current problems.

Analysts focus on processing the information stored in existing datasets and establishing the best way to present this data.

Data analysts rely on statistics and data modeling to solve problems and offer recommendations that can lead to immediate improvements.

Source: Google data analytics professional certificate

Data Science vs. Data Analytics

Data science

- Produces broad insights that concentrate on which questions should be asked about data
- Confronts what is unknown by using advanced techniques to make predictions about the future

Data analytics

- Emphasizes discovering answers to questions being asked
- Determines actionable insights that can be applied immediately based on existing queries

Source: Google data analytics professional certificate

The connections between data science and data analytics

Data science and data analytics share a fundamental goal: discover insights that can be used to lead an organization to improve and grow.

They are closely connected with information gathered through interactions within the measurable world.

As data projects become more complicated, organizations are discovering the advantages of assembling data teams, bringing data analysts and data scientists together.

Source: Google data analytics professional certificate

Data Science Pipeline



1. UNDERSTANDING
THE BUSINESS
GOALS



2. DATA GATHERING
AND EXPLORATION



3. FEATURE
ENGINEERING AND
MODEL SELECTION



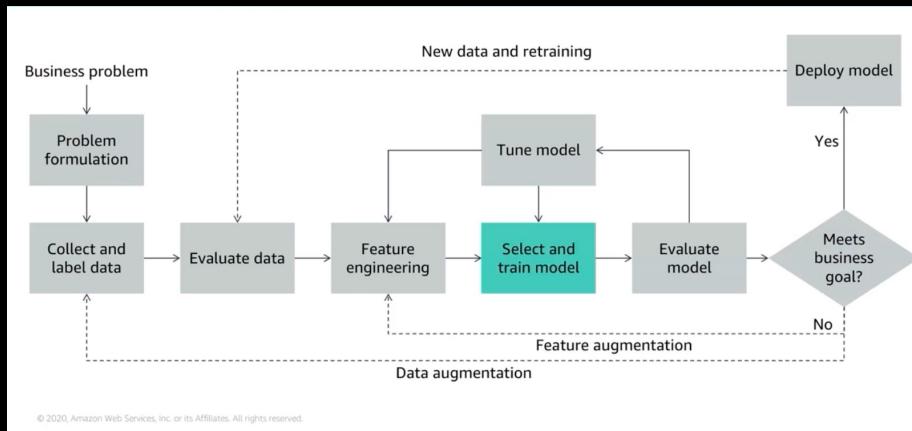
4. MODEL TRAINING,
EVALUATION, AND
REFINEMENT



5. MODEL
DEPLOYMENT AND
COMMUNICATION

Source: Google data analytics professional certificate

Machine Learning Pipeline



Dr. Vani Vasudevan

5/20/24 9

Source: AWS Academy – ML Foundation Course

Python For Data Science Cheat Sheet

Python Basics
Learn More Python for Data Science Interactively at www.datacamp.com

Variables and Data Types

Variable Assignment	>>> x=5	5
	>>> x	

Calculations With Variables

	x+2	Sum of two variables
	x-2	Subtraction of two variables
	x*2	Multiplication of two variables
	x**2	Exponentiation of a variable
	x%2	Remainder of a variable
	x/float(2)	Division of a variable
	2.5	

Types and Type Conversion

str()	'5', '3.45', 'True'	Variables to strings
int()	5, 3, 1	Variables to integers
float()	5.0, 1.0	Variables to floats
bool()	True, True, True	Variables to booleans

Asking For Help
>>> help(str)

Strings

```
>>> my_string = "thisStringIsAwesome"
>>> my_string
'thisStringIsAwesome'
```

String Operations

```
>>> my_string * 2
'thisStringIsAwesome' + 'thisStringIsAwesome'
>>> my_string + "Init"
'thisStringIsAwesomeInit'
>>> 'm' in my_string
True
```

Lists
[Also see NumPy Arrays](#)

Selecting List Elements	Index starts at 0
>>> a = 'is' >>> b = 'nice' >>> my_list = ['my', 'list', a, b] >>> my_list2 = [(4,5,6,7), [3,4,5,6]]	
	Select item at index 1 Select 3rd last item
>>> my_list[1:3]	Select items at index 1 and 2
>>> my_list[1:]	Select items after index 0
>>> my_list[:3]	Select items before index 3
>>> my_list[:]	Copy my_list
>>> my_list2[1][0]	my_list[[list][itemOfList]]
>>> my_list2[1][1:2]	

List Operations

>>> my_list + my_list	Get the index of an item
'[my', 'list', 'is', 'nice', 'my', 'list', 'is', 'nice']'	Count an item
>>> my_list * 2	Append an item at a time
'[my', 'list', 'is', 'nice', 'my', 'list', 'is', 'nice']'	Remove an item
>>> my_list > 4	Reverse the list
	Append an item
>>> my_list.pop(-1)	Remove an item
>>> my_list.insert('!')	Insert an item
>>> my_list.sort()	Sort the list

List Methods

>>> my_list.index('a')	Get the index of an item
>>> my_list.count('a')	Count an item
>>> my_list.append('!')	Append an item at a time
>>> my_list.remove('!')	Remove an item
>>> del(my_list[0:1])	Reverse the list
>>> my_list.reverse()	Append an item
>>> my_list.extend('!')	Remove an item
>>> my_list.pop(-1)	Insert an item
>>> my_list.insert(0, '!')	Sort the list

String Operations
[Also see Lists](#)

Selecting Numpy Array Elements	Index starts at 0
>>> my_list = [1, 2, 3, 4]	Get the index of an item
>>> my_array = np.array(my_list)	Count an item
>>> my_2darray = np.array([(1,2,3), [4,5,6]])	Select items at index 0 and 1

String Methods

>>> my_string.upper()	String to uppercase
>>> my_string.lower()	String to lowercase
>>> my_string.count('w')	Count String elements
>>> my_string.replace('e', 'i')	Replace String elements
>>> my_string.strip()	Strip whitespace

Libraries

Import Libraries	pandas Data analysis
>>> import numpy	Machine learning
>>> import numpy as np	
Selective Import	NumPy Scientific computing
>>> from math import pi	matplotlib 2D plotting

Install Python

ANACONDA Leading open data science platform powered by Python	spyder Free IDE that is included with Anaconda
	jupyter Create and share documents with live code, visualizations, text, ...

Numpy Arrays
[Also see Lists](#)

Selecting Numpy Array Elements	Index starts at 0
>>> my_array = np.array([1, 2, 3, 4])	Get the index of an item
>>> my_array[0:2]	Count an item
array([1, 2])	Select items at index 0 and 1
Subset 2D Numpy arrays	my_2darray[rows, columns]
>>> my_2darray[:, 0]	
array([1, 4])	

Numpy Array Operations

>>> my_array > 3	Get the dimensions of the array
array([False, False, False, True], dtype=bool)	Append items to an array
>>> my_array * 2	Insert items in an array
array([2, 4, 6, 8])	Delete items in an array
>>> my_array + np.array([5, 6, 7, 8])	Mean of the array
array([6, 8, 10, 12])	Median of the array
	Correlation coefficient
	Standard deviation

Numpy Array Functions

>>> my_array.shape	Get the dimensions of the array
>>> np.append(other_array)	Append items to an array
>>> np.insert(my_array, 1, 5)	Insert items in an array
>>> np.delete(my_array, [1])	Delete items in an array
>>> np.mean(my_array)	Mean of the array
>>> np.median(my_array)	Median of the array
>>> my_array.corrcoef()	Correlation coefficient
>>> np.std(my_array)	Standard deviation

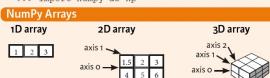
DataCamp
Learn Python for Data Science Interactively

Python For Data Science Cheat Sheet

Numpy Basics

Learn Python for Data Science Interactively at [www.DataCamp.com](#)

NumPy
 The NumPy library is the core library for scientific computing in Python. It provides a high-performance multidimensional array object, and tools for working with these arrays.
 Use the following import convention:
`>>> import numpy as np`

NumPy Arrays
 1D array 2D array 3D array


Creating Arrays

```
>>> a = np.array([1,2,3])
>>> b = np.array([(1,2,3), (4,5,6)], dtype = float)
>>> c = np.array([[1,2,3], [4,5,6]], [(1,2,3), (4,5,6)], dtype = float)
```

Initial Placeholders

```
>>> np.zeros((3,4)) Create an array of zeros
>>> np.ones((2,4),dtype=np.int16) Create an array of ones
>>> d = np.arange(10,15,2) Create an array of evenly spaced values (number of samples)
>>> e = np.linspace(0,2,9) Create an array of evenly spaced values (number of samples)
>>> f = np.full((2,2),7) Create a 2x2 identity matrix
>>> g = np.random.randint(2,2) Create an array with random values
>>> h = np.empty((3,2)) Create an empty array
```

I/O
Saving & Loading On Disk

```
>>> a = np.array([1,2,3])
>>> np.savetxt("my_array", a)
>>> np.load("myarray.npy")
```

Saving & Loading Text Files

```
>>> np.loadtxt("myfile.txt")
>>> np.genfromtxt("my_file.csv", delimiter=',')
>>> np.savetxt("myarray.txt", a, delimiter='=')
```

Data Types

```
>>> np.int64 Signed 64-bit integer type
>>> np.float32 Standard double-precision floating-point
>>> np.ndarray Container type represented by 1D arrays
>>> np.bool Boolean type storing TRUE and FALSE values
>>> np.object Python object type
>>> np.string_ Fixed-length string type
>>> np.unicode_ Fixed-length unicode type
```

Inspecting Your Array

>>> a.shape	Array dimensions
>>> len(a)	Length of array
>>> a.size	Number of array elements
>>> a.dtype	Data type of array elements
>>> a.itemsize	Native memory size
>>> b.astype(int)	Convert an array to a different type

Asking For Help

>>> np.info(np.ndarray.dtype)	
-------------------------------	--

Array Mathematics

>>> a = np.array([1,2,3]) >>> a + b	Addition
>>> a - b	Subtraction
>>> a * b	Multiplication
>>> a / b	Division
>>> a % b	Modulo
>>> a ** b	Pow
>>> a * np.sqrt(b)	Square root
>>> a * np.sin(a)	Element-wise sine
>>> a * np.log(a)	Element-wise natural logarithm
>>> a * np.exp(b)	Dot product

Arithmetic Operations

>>> a = np.array([-3,-1,0,1,3]) >>> np.subtract(a,b)	-
>>> a + b	+
>>> a * b	*
>>> a / b	/
>>> a % b	%
>>> a ** b	**
>>> a * np.sqrt(b)	sqrt
>>> a * np.sin(a)	sin
>>> a * np.log(a)	log
>>> a * np.exp(b)	exp

Comparison

>>> a == b	Element-wise comparison
>>> a != b	Not equal
>>> a < b	Less than
>>> a > b	Greater than
>>> a <= b	Less than or equal
>>> a >= b	Greater than or equal

Aggregating Functions

>>> a.sum()	Aggregate sum
>>> a.min()	Aggregate minimum
>>> a.max(axis=0)	Maximum value of an array row
>>> a.cumsum(axis=1)	Cumulative sum of the elements
>>> a.mean()	Mean
>>> a.std()	Standard deviation
>>> a.var()	Variance
>>> a.corcoef()	Correlation coefficient
>>> a.std(ddof=1)	Standard deviation (ddof=1)

Copying Arrays

>>> b = a.view() >>> b.copy() >>> b = a.copy()	Create a view of the array with the same data Create a copy of the array Create a deep copy of the array
--	--

Sorting Arrays

>>> a.argsort()	Sort an array
>>> a.argsort(axis=0)	Sort the elements of an array's axis

Subsetting, Slicing, Indexing

>>> a[2]	Select the element at the 2nd index
>>> a[1,2]	Select the element at row 1 column 2 (equivalent to b[1][2])
>>> a[0]	Select items at index 0 and 1 in column 1
>>> b[1,2,1]	Select all items at row 0 (equivalent to a[[0], :])
>>> b[1,1]	Same as [[1], [1]]
>>> b[1:,1]	Reversed array a
>>> a[a < 2]	Select elements from a less than 2
>>> a[[1,0,0,1,1,0,0]]	Select elements (1,0,0,0,1,1,0,0) and (0,0,0)

Also see Lists

Array Manipulation

>>> np.transpose(b)	Permute array dimensions
>>> b.T	Permute array dimensions
>>> g.reshape(3,-2)	Flatten the array
>>> h.resize((2,6))	Reshape, but don't change data
>>> i.insert(0,1,5)	Return a new array with shape (2,6)
>>> j.delete(0)	Append items to an array
>>> k.concatenate((a,d))	Insert items in an array
>>> l.concatenate((a,[1]))	Delete items from an array
>>> m.concatenate((a,b))	Concatenate arrays
>>> n.vstack((a,b))	Stack arrays vertically (row-wise)
>>> o.hstack((a,b))	Stack arrays vertically (row-wise)
>>> p.column_stack((a,d))	Stack arrays horizontally (column-wise)
>>> q.column_stack((a,d))	Stack arrays horizontally (column-wise)
>>> r.stack((a,d))	Create stacked column-wise arrays
>>> s.stack((a,d))	Create stacked column-wise arrays
>>> t.hsplit(a,3)	Split the array horizontally at the 3rd index
>>> u.vsplit(a,3)	Split the array vertically at the 2nd index

DataCamp
 Learn Python for Data Science Interactively

Python For Data Science Cheat Sheet

Pandas Basics

Learn Python for Data Science interactively at www.DataCamp.com

Pandas

The Pandas library is built on NumPy and provides easy-to-use data structures and data analysis tools for the Python programming language.

Pandas Data Structures

Series

A one-dimensional labeled array capable of holding any data type

Index

```
>>> s = pd.Series([3, -5, 7, 4], index=['a', 'b', 'c', 'd'])
```

DataFrame

Columns

Country	Capital	Population
Belgium	Brussels	11000000
India	New Delhi	1303171035
Brazil	Brasilia	207847528

Index

```
>>> df = pd.DataFrame(data, columns=['Country', 'Capital', 'Population'])
```

I/O

Read and Write to CSV

```
>>> pd.read_csv('file.csv', header=None, nrows=5)
>>> df.to_csv('myDataFrame.csv')
```

Read and Write to Excel

```
>>> pd.read_excel('file.xlsx')
>>> pd.to_excel('file/myDataFrame.xlsx', sheet_name='Sheet1')
>>> Read multiple sheets from the same file
>>> xlxs = pd.ExcelFile('file.xls')
>>> df = pd.read_excel(xlxs, 'Sheet1')
```

Asking For Help

```
>>> help(pd.Series.loc)
```

Selection

Getting

```
>>> s['b']
>>> df[1]
>>> df[[1]]
```

Get one element

Get subset of a DataFrame

Selecting, Boolean Indexing & Setting

By Position

```
>>> df.iloc[[0], [0]]
>>> df.iat[[0], [0]]
```

Select single value by row & column

By Label

```
>>> df.loc[[0], ['Country']]
>>> df.loc[[0], 'Belgium']
>>> df.at[[0], 'Country']
>>> df['Belgium']
```

Select single value by row & column labels

By Label/Position

```
>>> df.ix[2]
>>> df.ix[0, 'Country']
>>> df.ix[0, 'Belgium']
>>> df.ix[[0], ['Country']]
>>> df.ix[0, 'Belgium']
```

Select single row of subset of rows

Boolean Indexing

```
>>> s[(s > 1)]
>>> s[(s < -1) | (s > 2)]
>>> df[df['Population']>1200000000]
```

Series s where value is not > 1 or < -1 or > 2. Use filter to adjust DataFrame

Setting

```
>>> s['a'] = 6
```

Set index a of Series s to 6

Dropping

```
>>> s.drop(['a', 'c'])
>>> df.drop('Country', axis=1)
```

Drop values from rows (axis=0) Drop values from columns (axis=1)

Sort & Rank

```
>>> df.sort_index()
>>> df.sort_values(by='Country')
>>> df.rank()
```

Sort by labels along an axis Sort by the values along an axis Assign ranks to entries

Retrieving Series/DataFrame Information

Basic Information

```
>>> df.shape
>>> df.index
>>> df.columns
>>> df.info()
>>> df.count()
```

(rows,columns) Describe index Describe columns Info on DataFrame Number of non-NA values

Summary

```
>>> df.sum()
>>> df.cumsum()
>>> df.min() / df.max()
>>> df.idxmin() / df.idxmax()
>>> df.describe()
>>> df.quantile()
>>> df.median()
```

Sum of values Cumulative sum of values Minimum/maximum value Minimum/maximum index value Summary statistics Quantiles Median of values

Applying Functions

```
>>> f = lambda x: x*x
>>> df.apply(f)
>>> df.applymap(f)
```

Apply function Apply function element-wise

Data Alignment

Internal Data Alignment

NA values are introduced in the indices that don't overlap:

```
>>> s3 = pd.Series([7, -2, 3], index=['a', 'c', 'd'])
>>> s = s3
>>> a = 10.0
>>> b = 5.0
>>> c = 5.0
>>> d = 7.0
```

Arithmetic Operations with Fill Methods

You can also do the internal data alignment yourself with the help of the fill methods:

```
>>> s.add(s3, fill_value=0)
>>> s + s3
>>> b = -5.0
>>> c = 5.0
>>> d = 7.0
>>> s.sub(s3, fill_value=2)
>>> s.div(s3, fill_value=4)
>>> s.mul(s3, fill_value=3)
```

DataCamp
Learn Python for Data Science interactively

Python For Data Science Cheat Sheet

Matplotlib

Learn Python interactively at www.DataCamp.com



Matplotlib

Matplotlib is a Python 2D plotting library which produces publication-quality figures in a variety of hardcopy formats and interactive environments across platforms.

1 Prepare The Data

1D Data

```
>>> import numpy as np
>>> x = np.linspace(0, 100)
>>> y = np.sqrt(x)
```

3D Data or Images

```
>>> data = np.random.rand(10, 10)
>>> data = 3 + np.random.rand(10, 10)
>>> data = np.reshape(data, (10, 10))
>>> U = 1 - X**2 - Y**2
>>> from matplotlib.colors import get_sample_data
>>> img = np.load(get_sample_data("axes_grid/bivariate_normal.npy"))
```

2 Create Plot

Figure

```
>>> import matplotlib.pyplot as plt
```

Figure

```
>>> fig = plt.figure()
>>> fig = plt.figure(figsize=plt.rcParams["figure.figsize"])
```

Axes

All plotting is done with respect to an `Axes`. In most cases, a subplot will fit your needs. A subplot is an axes on a grid system.

```
>>> fig.add_axes()
>>> ax = fig.add_subplot(221) # row-col-num
>>> ax.set_xlabel('X')
>>> ax.set_ylabel('Y')
>>> fig1, axes = plt.subplots(nrows=2, ncols=2)
>>> fig2, axes2 = plt.subplots(ncols=2)
```

3 Plotting Routines

1D Data

```
>>> x = np.linspace(0, 10)
>>> lines = ax.plot(x)
>>> ax.plot([1, 2, 3], [4, 5, 6])
>>> ax.fill_between(x, 0, 10, color='red', alpha=0.5)
>>> ax.fill_between(x, 0, 10, color='blue', alpha=0.5)
>>> ax.fill_between(x, 0, 10, color='yellow')
```

2D Data or Images

```
>>> fig, ax = plt.subplots()
>>> ax.imshow(img,
             cmap='gist_earth',
             interpolation='nearest',
             vmin=0,
             vmax=1)
```

Plot Anatomy & Workflow

Plot Anatomy



Figure

Workflow

The basic steps to creating plots with matplotlib are:

- 1 Prepare data
- 2 Create plot
- 3 Plot
- 4 Customize plot
- 5 Save plot
- 6 Show plot

```
>>> import matplotlib.pyplot as plt
>>> x = [1, 2, 3, 4, 5]
>>> y = [10, 20, 25, 30, 35]
>>> fig = plt.figure()
>>> ax = fig.add_subplot(111)
>>> ax.plot(x, y, 'r--', linewidth=3)
>>> ax.set_title('sigma=15$')
>>> ax.set_xlabel('X-axis')
>>> ax.set_ylabel('Y-axis')
>>> plt.show()
```

4 Customize Plot

Colors, Color Bars & Color Maps

```
>>> plt.plot(x, X, X*x, X*x**2, X*x**3)
>>> plt.plot(X, y, 'r', c='red')
>>> ax.plot(X, y, 'c', c='cyan')
>>> ax.plot(X, y, 'k', c='black')
>>> im = ax.imshow(data,
                  cmap='seismic')
```

Markers

```
>>> fig = plt.subplot()
>>> ax.scatter(x,y,marker="*")
>>> ax.plot(x,y,marker="o")
```

Linestyles

```
>>> plt.plot(x,y,linewidth=4.0)
>>> plt.plot(x,y,linestyle="dashed")
>>> plt.plot(x,y,linestyle="dashdot")
>>> plt.plot(x,y,linestyle="solid")
>>> plt.step(x,y,color="r",linewidth=4.0)
```

Text & Annotations

```
>>> ax.text(2,1,
          "Create graph",
          style="italic")
>>> ax.annotate("text",
              xy=(3, 0),
              xytext=(0, 1),
              textcoords="offset",
              arrowprops=dict(arrowstyle=">",
                             connectionstyle="arc3"))
```

Vector Fields

```
>>> axes[0,1].arrow(0,0,0.5,0.5)
>>> axes[0,1].quiver(x,y)
```

Distribution

```
>>> axes[0,1].streamplot(X,Y,U,V)
>>> axes[1,0].hist(y)
>>> axes[1,0].violinplot(z)
```

5 Save Plot

Save Figures

```
>>> plt.savefig('foo.png')
```

Save transparent figures

```
>>> plt.savefig('foo.png', transparent=True)
```

6 Show Plot

Close & Clear

Clear on exit

```
>>> plt.close()
>>> plt.close('all')
>>> plt.close()
```

Clear the entire figure

Close a window

Python For Data Science Cheat Sheet

Seaborn

Learn Data Science Interactively at www.DataCamp.com



Statistical Data Visualization With Seaborn

The Python visualization library Seaborn is based on matplotlib and provides a high-level interface for drawing attractive statistical graphics.

Make use of the following aliases to import the libraries:

```
>>> import matplotlib.pyplot as plt
>>> import seaborn as sns
```

The basic steps to creating plots with Seaborn are:

1. Prepare some data
2. Control figure aesthetics
3. Plot with Seaborn
4. Further customize your plot

Importing Seaborn:

```
>>> import matplotlib.pyplot as plt
>>> tips = sns.load_dataset("tips")  
Step 1
>>> sns.set_style("whitegrid")  
Step 2
>>> g = sns.FacetGrid(tips, col="smoker", row="sex")
>>> g.map(plt.hist, "total_bill", bins=10)  
Step 3
>>> g = g.map(plt.bar, "size", "tip", aspect=1)
>>> g.set(xlim=(0,10), ylim=(0,100))  
Step 4
>>> plt.title("title")  
Step 5
>>> g.show()
```

1) Data

Also see [Lists](#), [NumPy](#) & [Pandas](#)

```
>>> import pandas as pd
>>> import numpy as np
>>> uniform_data = np.random.rand(10, 12)
>>> data = pd.DataFrame(np.arange(1,101),
>>>                      columns=[str(i) for i in range(1,10)])
>>> g = sns.FacetGrid(data, hue="sex", palette="Greens_d")
>>> g.map(plt.hist, "age", bins=10)
>>> g.map(plt.violinplot, "age", "survived", inner="quartile",
>>>       data=titanic)
>>> g.map(plt.violinplot, "age", "survived", inner="box",
>>>       data=titanic)
```

Seaborn also offers built-in data sets:

```
>>> titanic = sns.load_dataset("titanic")
>>> tips = sns.load_dataset("tips")
```

2) Figure Aesthetics

Also see [Matplotlib](#)

```
>>> f, ax = plt.subplots(figsize=(5, 6))  
Create a figure and one subplot
Seaborn styles
>>> sns.set()  
(Reset the seaborn default
>>> sns.set_style("whitegrid")  
(Set the matplotlib parameters
>>> sns.set_style(style, **kwargs)  
(Set the matplotlib parameters
>>> sns.set(**kwargs)  
(Return a dict of params or use with
>>> sns.axes_style("whitegrid")  
with to temporarily set the style
```

3) Plotting With Seaborn

Axis Grids

<pre>>>> g = sns.FacetGrid(titanic, >>> col="survived", >>> row="sex")</pre>	Subplot grid for plotting conditional relationships
<pre>>>> g = g.map(plt.hist, "petal_length", >>> hue="survived", >>> data=titanic)</pre>	Draw a categorical plot onto a FacetGrid
<pre>>>> sns.lmplot(x="equal_width", >>> y="equal_length", >>> hue="species", >>> data=titanic)</pre>	Plot data and regression model fits across a FacetGrid
<pre>>>> b = sns.PairGrid(titanic) >>> h = b.map(plt.scatter)</pre>	Subplot grid for plotting pairwise relationships
<pre>>>> i = sns.JointGrid(*iris) >>> i = i.map(plt.kdeplot,</pre>	Plot pairwise bivariate distributions
<pre> data=titanic, y="y", data=titanic)</pre>	Grid for bivariate plot with marginal univariate plots
<pre>>>> i = i.map(sns.residplot,</pre>	
<pre> sns.jointplot("equal_length", "equal_width", data=titanic, kind="hex")</pre>	Plot bivariate distribution

Categorical Plots

<pre>>>> sns.stripplot(x="species", >>> y="petal_length", >>> data=titanic)</pre>	Scatterplot with one categorical variable
<pre>>>> sns.swarmplot(x="species", >>> y="petal_length", >>> data=titanic)</pre>	Categorical scatterplot with non-overlapping points
<pre>>>> sns.factorplot(x="class", >>> y="survived", >>> hue="class", >>> data=titanic)</pre>	Show point estimates and confidence intervals with scatterplot glyphs
<pre>>>> sns.countplot(x="deck", >>> data=titanic, >>> palette="Greens_d")</pre>	Show count of observations
<pre>>>> sns.pointplot(x="class", >>> y="survived", >>> hue="sex", >>> data=titanic, >>> palette=({"male": "#e74c3c", "female": "#3182bd"}, {"survived": "#f7a5c2", "not survived": "#b159bd"}, {"deck": "#3182bd"}, {"age": "#3182bd"}, {"marker": "#3182bd"}, {"linestyle": ["--", "-"]})</pre>	Show point estimates and confidence intervals as rectangular bars
<pre>>>> sns.bxpplot(x="alive", >>> y="age", >>> hue="alive", >>> data=titanic)</pre>	Bxpplot
<pre>>>> sns.bxpplot(x="alive", >>> y="age", >>> hue="alive", >>> data=titanic)</pre>	Bxpplot with wide-form data
<pre>>>> sns.violinplot(x="sex", >>> y="survived", >>> data=titanic)</pre>	Violin plot

Regression Plots

<pre>>>> sns.regression(x="equal_width", >>> y="equal_length", >>> data=iris, >>> ax=ax)</pre>	Plot data and a linear regression model fit
---	---

Distribution Plots

<pre>>>> plot = sns.distplot(data.y, >>> kde=True, >>> color="blue")</pre>	Plot univariate distribution
---	------------------------------

Matrix Plots

<pre>>>> sns.heatmap(uniform_data, vmin=0, vmax=1)</pre>	Heatmap
---	---------

4) Further Customizations

Also see [Matplotlib](#)

Axis Grid Objects

<pre>>>> g.despine(left=True) >>> g.set_ylabels("Survived") >>> g.set_xticklabels(rotation=45) >>> g.set_axis_labels("Deck") >>> h.set_xlim(0,5) >>> h.set_ylim(0,5) >>> h.set_xticks([0,2.5,5]) >>> h.set_yticks([0,2.5,5])</pre>	Remove left spine Set the labels of the y-axis Set the tick labels for x-axis Set the axis labels for both Set the limits and ticks of the x-and-y-axes
--	---

Plot

<pre>>>> plt.title("Title") >>> plt.xlabel("X-axis label") >>> plt.ylabel("Y-axis label") >>> plt.ylim(0,100) >>> plt.xlim(0,100) >>> plt.setp(ax,yticks=[0,5]) >>> plt.tight_layout()</pre>	Add plot title Adjust the label of the y-axis Adjust the label of the x-axis Adjust the limits of the x-axis Adjust the limits of the y-axis Adjust a plot property Adjust subplot params
---	---

5) Show or Save Plot

Also see [Matplotlib](#)

<pre>>>> plt.show() >>> plt.savefig("foo.png") >>> plt.savefig("foo.png", >>> transparent=True)</pre>	Show the plot Save the plot as a figure Save transparent figure
---	---

Close & Clear

Also see [Matplotlib](#)

<pre>>>> plt.close() >>> plt.cla() >>> plt.close()</pre>	Clear an axis Clear the current figure Close a window
---	---

DataCamp
Learn Python for Data Science Interactively!

Python For Data Science Cheat Sheet

3) Renderers & Visual Customizations

Glyphs

```
>>> p1.circle("mpg", "cyl", source=cds_df,
    factors=[87], color="blue", size=1)
>>> p2.square("mpg", "cyl", source=cds_df,
    color="blue", size=1)
```

Line Glyphs

```
>>> p1.line([(1,2,3,4), (3,4,5,6)], line_width=2)
>>> p2.multi_line([(1,2,3,7), (5,6,7)]),
    pd.DataFrame([(3,4,5),(3,5,6)]),
    color="blue")
```

Customized Glyphs

Also see Data

```
>>> p = figure(tools='box_select')
>>> p.circle("mpg", "cyl", source=cds_df,
    selection_color='red',
    nonselection_alpha=0.1)
```

Hover Glyphs

```
>>> from bokeh.models import HoverTool
>>> hover = HoverTool(tooltips=None, mode='vline")
```

ColorMapping

```
>>> from bokeh.models import CategoricalColorMapper
>>> color_mapper = CategoricalColorMapper(factors=[87], palette=['blue', 'red', 'green'])
>>> p3.circle("mpg", "cyl", source=cds_df,
    color_dict={'field': 'origin',
    'transform': color_mapper},
    legend='Origin')
```

Legend Location

Inside Plot Area

```
>>> p.legend.location = "bottom_left"
```

Outside Plot Area

```
>>> from bokeh.models import Legend
>>> r1 = pd.line([(1,2,3,4), [3,4,5,6)])
>>> r2 = pd.line([(1,2,3,7), (5,6,7)])
>>> legend = Legend(items=[("r1", r1), ("r2", r2)], location=(0, -50))
>>> p.add_layout(legend, "right")
```

Legend Orientation

```
>>> p.legend.orientation = "horizontal"
>>> p.legend.orientation = "vertical"
```

Legend Background & Border

```
>>> p.legend.border_line_color = "navy"
>>> p.legend.background_fill_color = "white"
```

Rows & Columns Layout

Rows

```
>>> from bokeh.layouts import row
>>> layout = row(p1,p2,p3)
```

Column

```
>>> from bokeh.layouts import column
>>> layout = column(p1,p2,p3)
```

Nesting Rows & Columns

```
>>> layout = row(column(p1,p2), p3)
```

Grid Layout

```
>>> from bokeh.layouts import gridplot
>>> row1 = [p1,p2]
>>> row2 = [p3]
>>> layout = gridplot([[p1,p2],[p3]])
```

Tabbed Layout

```
>>> from bokeh.models.widgets import Panel, Tabs
>>> tab1 = Panel(child=p1, title="tab1")
>>> tab2 = Panel(child=p2, title="tab2")
>>> layout = Tabs(tabs=[tab1, tab2])
```

Linked Plots

Linked Axes

```
>>> p2.x_range = p1.x_range
>>> p2.y_range = p1.y_range
```

Linked Brushing

```
>>> p4 = Figure(plot_width = 100,
    tools='box_select,lasso_select')
>>> p4.circle("mpg", "cyl", source=cds_df)
>>> p5 = Figure(plot_width = 200,
    tools='box_select,lasso_select')
>>> p5.circle("mpg", "cyl", source=cds_df)
>>> layout = row(p4,p5)
```

4) Output & Export

Notebook

```
>>> from bokeh.io import output_notebook, show
>>> output_notebook()
```

HTML

Standalone HTML

```
>>> from bokeh.embed import file_html
>>> from bokeh.resources import CDN
>>> html = file_html(p, CDN, "my_plot")
```

PNG

```
>>> from bokeh.io import export_png
>>> export_png(p, filename="plot.png")
```

SVG

```
>>> from bokeh.io import export_svgs
>>> p.output_backend = "svg"
>>> export_svgs(p, filename="plot.svg")
```

5) Show or Save Your Plots

```
>>> show(p1) >>> show(layout)
>>> save(p1) >>> save(layout)
```

DataCamp
Learn Python for Data Science Interactively

Python For Data Science Cheat Sheet

SciPy - Linear Algebra

Learn More Python for Data Science [Interactively at www.datacamp.com](#)



SciPy
The SciPy library is one of the core packages for scientific computing that provides mathematical algorithms and convenience functions built on the NumPy extension of Python.

Interacting With NumPy [\(Also see NumPy\)](#)

```
>>> import numpy as np
>>> a = np.array([1,2,3])
>>> b = np.array([(4,5,6),(2,3)])
>>> c = np.array([(1,2,3),(4,5,6),[(3,2,1), (4,5,6)]])
```

Index Tricks

```
>>> np.mgrid[0:5,0:5] Create a dense meshgrid
>>> a = np.arange(10) Create an open meshgrid
>>> np.r_[[1,0]*5,-1:1:10j] Stack arrays vertically (row-wise)
>>> np.c_[b,c] Create stacked column-wise arrays
```

Shape Manipulation

```
>>> np.transpose(b) Permute array dimensions
>>> b.flatten() Flatten the array
>>> np.vstack((b,c)) Stack arrays horizontally (column-wise)
>>> np.vstack(c) Stack arrays vertically (row-wise)
>>> np.hsplit(c,2) Split the array horizontally at the 2nd index
>>> np.vsplit(c,2) Split the array vertically at the 2nd index
```

Polynomials

```
>>> from numpy import polyd Create a polynomial object
>>> p = polyd([1,4,5])
```

Vectorizing Functions

```
>>> def myfunc(a):
...     if a < 0:
...         return a**2
...     else:
...         return a+2
>>> np.vectorize(myfunc) Vectorize functions
```

Type Handling

```
>>> np.real(z) Return the real part of the array elements
>>> np.imag(z) Return the imaginary part of the array elements
>>> np.real_if_close(a, tol=1000) Return a real array if complex parts close to 0
>>> np.cast["f"](np.pi) Cast object to a data type
```

Other Useful Functions

```
>>> np.angle(b,deg=True) Return the angle of the complex argument
>>> g = np.linspace(0,np.pi,num=5) Create an array of evenly spaced values
>>> np.unwrap(g) Unwrap
>>> np.logspace(0,10,3) Create an array of evenly spaced values (log scale)
>>> np.random.choice(a,[1,2]) Returns values from a list of arrays depending on conditions
>>> misc.factorial(5) Factorial
>>> misc.comb(10,3,exact=True) Compute N things taken at k time
>>> misc.central_diff_weights(3) Weights for N-point central derivative
>>> np.derivative(myfunc,1,0) Find the n-th derivative of a function at a point
```

Linear Algebra
You'll use the `linalg` and `sparse` modules. Note that `scipy.linalg` contains and expands on `numpy.linalg`.

Creating Matrices

```
>>> A = np.matrix(np.random.random((2,2)))
>>> B = np.amsmatrix(B)
>>> C = np.mat(np.random.random((10,5)))
>>> D = np.mat((3,4), [5,6]))
```

Basic Matrix Routines

	Matrix Functions	
Inverse	<code>np.linalg.inv(A)</code>	Addition
Inverse	<code>np.linalg.inv(A)</code>	Subtraction
Transpose matrix	<code>np.linalg.transpose(A)</code>	Division
Conjugate transposition	<code>np.linalg.conjugate_transpose(A)</code>	Multiplication
Trace	<code>np.trace(A)</code>	Dot product

Multiplication

```
>>> np.multiply(D,A)
>>> np.dot(A,D)
>>> np.vdot(A,D)
>>> np.inner(A,D)
>>> np.outer(A,D)
>>> np.tensordot(A,D)
>>> np.kron(A,B)
```

Exponential Functions

```
>>> np.linalg.expm(A)
>>> np.linalg.expm(A)
>>> np.linalg.expm1(A)
>>> np.linalg.expm2(A)
```

Logarithmic Function

```
>>> np.linalg.logm(A)
```

Trigonometric Functions

```
>>> np.linalg.sinm(D)
>>> np.linalg.cosm(D)
>>> np.linalg.tanm(A)
>>> np.linalg.acosm(D)
>>> np.linalg.asinm(D)
>>> np.linalg.atanm(A)
```

Hyperbolic Trigonometric Functions

```
>>> np.linalg.sinh(D)
>>> np.linalg.cosh(D)
>>> np.linalg.tanhm(A)
```

Matrix Sign Function

```
>>> np.linalg.signm(A)
```

Matrix Square Root

```
>>> np.linalg.sqrtm(A)
```

Arbitrary Functions

```
>>> np.linalg.funm(A, lambda x: x*x)
```

Decompositions

Eigenvalues and Eigenvectors	<code>la.eig(A)</code>	Solve ordinary or generalized eigenvalue problem for square matrix
Inverse	<code>np.linalg.inv(I)</code>	Unpack eigenvectors
Norm	<code>np.linalg.norm(I)</code>	First eigenvector
Solving linear problems	<code>np.linalg.solve(H,x)</code>	Second eigenvector
		Unpack eigenvectors
Sparse Matrix Routines	<code>np.linalg.linsolve.spsolve(H,x)</code>	
Inverse	<code>np.linalg.inv(C)</code>	
Norm	<code>np.linalg.norm(C)</code>	

Solving linear problems

```
>>> I = sparse.csr_matrix(D)
>>> H = sparse.csr_matrix(D)
>>> E = sparse.dok_matrix(D)
>>> E.todense()
>>> np.linalg.linsolve.spsolve(H,x)
```

Compressed Sparse Row matrix

Compressed Sparse Column matrix

Dok matrix

Solver for full matrix

Solver for sparse matrix

Sparse Matrix Functions

	<code>np.linalg.expm(I)</code>	Sparse matrix exponential

Asking For Help

[DataCamp](#)
Learn Python for Data Science [Interactively](#)

[www.datacamp.com](#)

Outline



Setting the Context



Introduction to advanced Python libraries



High-performance data manipulation with Vaex



Interactive data exploration with D-Tale



Feature engineering with Feature Engine

Importance of Advanced Python Libraries in Enhancing Data Analysis Capabilities and Improving Workflow Efficiency

1. Efficiency and Performance:

•Handling Large Datasets:

- Advanced libraries like Vaex are designed to work with large datasets that do not fit into memory, enabling efficient out-of-core data manipulation and analysis. This capability is crucial for processing large-scale data in a timely manner.

•Optimized Computations:

- These libraries leverage optimized algorithms and data structures, significantly reducing computation time. For example, Vaex uses memory-mapped files and efficient indexing to speed up data operations.

•Parallel Processing:

- Many advanced libraries support parallel processing and distributed computing, allowing for faster data processing by utilizing multiple CPU cores or even clusters of machines.

2. Ease of Use and Flexibility:

•User-Friendly APIs:

- Libraries like D-Tale provide intuitive interfaces for data exploration and

visualization, making it easier for data professionals to interact with their data without extensive coding.

•Seamless Integration:

- These libraries integrate seamlessly with other popular Python libraries such as Pandas, NumPy, and scikit-learn, allowing for a smooth transition and interoperability between different tools in the data science ecosystem.

3. Community Support and Documentation:

•Comprehensive Documentation:

- Most advanced libraries come with extensive documentation and tutorials, providing clear guidelines on how to use their features effectively. This makes it easier for users to get started and overcome any challenges they may encounter.

•Active Community:

- An active community of users and developers contributes to continuous improvement and support. Users can seek help from community forums, GitHub issues, and other online resources, ensuring they are not alone in their learning journey.

4. Improving Workflow Efficiency:

•Streamlined Data Processing:

- Feature Engine simplifies common preprocessing tasks such as handling missing data, encoding categorical variables, and scaling features, making the data preparation phase more efficient and less error-prone.

•Interactive Analysis:

- D-Tale provides an interactive environment for data exploration, enabling quick insights and facilitating a more iterative and exploratory analysis process. This leads to faster discovery of patterns and anomalies in the data.

5. Enhancing Analytical Capabilities:

•Advanced Techniques:

- By leveraging advanced libraries, data professionals can implement sophisticated techniques and models that were previously challenging or time-consuming. This includes out-of-core algorithms, automated machine learning, and complex feature engineering.

•Better Decision Making:

- The ability to analyze large datasets efficiently and accurately leads to more informed and data-driven decisions. This is particularly important in industries like healthcare, where timely and accurate insights can significantly impact business outcomes.

•Scalability:

- Advanced libraries are designed to scale with the data. As the volume and

complexity of data grow, these tools ensure that data analysis remains manageable and performant.

Introduction to Advanced Libraries and Frameworks⁽¹⁾

Vaex

19

A high-performance Data Frame library designed for out-of-core data manipulation. Vaex is optimized for working with large datasets that do not fit into memory, making it an excellent choice for big data projects.

[repo](#)

```
conda create -n myenv python=3.10
conda activate myenv
```

```
pip install pydantic-settings
pip show pydantic-settings
conda install -c conda-forge vaex
vaex --version
```

<https://vaex.readthedocs.io/en/latest/datasets.html>

Introduction to Advanced Libraries and Frameworks⁽²⁾

D-Tale

20

An interactive web-based tool for visualizing and analysing Pandas Data Frames. D-Tale provides a user-friendly interface that allows you to explore data interactively, making the data exploration process more efficient and insightful.

repo

<https://www.nyc.gov/site/tlc/about/tlc-trip-record-data.page>

Introduction to Advanced Libraries and Frameworks⁽³⁾

Feature Engine

21

A library for feature engineering that integrates seamlessly with scikit-learn. Feature Engine offers various transformers for handling missing data, encoding categorical variables, and scaling features, simplifying the data preprocessing phase.

[repo](#)

Outline



Setting the Context



Introduction to
advanced Python
libraries



High-performance
data manipulation
with Vaex



Interactive data
exploration with D-
Tale



Feature engineering
with Feature Engine

Importance of Advanced Python Libraries in Enhancing Data Analysis Capabilities and Improving Workflow Efficiency

1. Efficiency and Performance:

•Handling Large Datasets:

- Advanced libraries like Vaex are designed to work with large datasets that do not fit into memory, enabling efficient out-of-core data manipulation and analysis. This capability is crucial for processing large-scale data in a timely manner.

•Optimized Computations:

- These libraries leverage optimized algorithms and data structures, significantly reducing computation time. For example, Vaex uses memory-mapped files and efficient indexing to speed up data operations.

•Parallel Processing:

- Many advanced libraries support parallel processing and distributed computing, allowing for faster data processing by utilizing multiple CPU cores or even clusters of machines.

2. Ease of Use and Flexibility:

•User-Friendly APIs:

- Libraries like D-Tale provide intuitive interfaces for data exploration and

visualization, making it easier for data professionals to interact with their data without extensive coding.

•Seamless Integration:

- These libraries integrate seamlessly with other popular Python libraries such as Pandas, NumPy, and scikit-learn, allowing for a smooth transition and interoperability between different tools in the data science ecosystem.

3. Community Support and Documentation:

•Comprehensive Documentation:

- Most advanced libraries come with extensive documentation and tutorials, providing clear guidelines on how to use their features effectively. This makes it easier for users to get started and overcome any challenges they may encounter.

•Active Community:

- An active community of users and developers contributes to continuous improvement and support. Users can seek help from community forums, GitHub issues, and other online resources, ensuring they are not alone in their learning journey.

4. Improving Workflow Efficiency:

•Streamlined Data Processing:

- Feature Engine simplifies common preprocessing tasks such as handling missing data, encoding categorical variables, and scaling features, making the data preparation phase more efficient and less error-prone.

•Interactive Analysis:

- D-Tale provides an interactive environment for data exploration, enabling quick insights and facilitating a more iterative and exploratory analysis process. This leads to faster discovery of patterns and anomalies in the data.

5. Enhancing Analytical Capabilities:

•Advanced Techniques:

- By leveraging advanced libraries, data professionals can implement sophisticated techniques and models that were previously challenging or time-consuming. This includes out-of-core algorithms, automated machine learning, and complex feature engineering.

•Better Decision Making:

- The ability to analyze large datasets efficiently and accurately leads to more informed and data-driven decisions. This is particularly important in industries like healthcare, where timely and accurate insights can significantly impact business outcomes.

•Scalability:

- Advanced libraries are designed to scale with the data. As the volume and

complexity of data grow, these tools ensure that data analysis remains manageable and performant.

Key Features of Vaex

Feature	Description
Out-of-Core Processing	Vaex allows you to work with datasets larger than your system's RAM by using out-of-core processing, reading and processing data in chunks.
Memory-Mapped Files	Utilizes memory-mapped files for fast and efficient data access, speeding up data loading and processing.
Lazy Evaluation	Operations in Vaex are executed lazily, meaning computations are performed only when results are needed, optimizing execution and reducing memory usage.
No Copying of Data	Avoids unnecessary data copying, enhancing performance and efficiency during operations on large datasets.
Data Exploration and Visualization	Integrates with visualization libraries like Matplotlib, Plotly, and Bokeh, providing powerful tools for interactive data visualization and exploration.
Advanced Filtering and Selection	Supports complex filtering and selection operations, including boolean indexing, range selection, and conditional filtering, optimized for performance.
Fast Aggregations and GroupBy Operations	Provides fast and efficient aggregation functions and groupby operations for quick summarization and analysis of large datasets.
Virtual Columns	additional memory, useful for performing calculations and transformations on existing data.
Integration with Pandas	Interoperates with Pandas, allowing seamless conversion between Vaex DataFrames and Pandas DataFrames when needed.
Support for Multiple Data Formats	HDF5, FITS, Parquet, and Arrow, providing flexibility in data storage and access.
Parallel Processing	Leverages multiple cores and processors for parallel computation, enhancing its ability to handle large datasets efficiently.

5/20/24

Dr. Vani Vasudevan

23

Sources:

<https://vaex.readthedocs.io/en/latest/>

<https://vaex.readthedocs.io/en/latest/guides/performance.html>

<https://vaex.io/blog/8-incredibly-powerful-Vaex-features-you-might-have-not-known-about>

Key differences between Pandas and Vaex

Feature	Pandas	Vaex
Data Handling	In-memory	Out-of-core
Dataset Size	Limited by RAM	Handles datasets larger than RAM
Performance	Slower for very large datasets	Optimized for large datasets
Computation	Eager evaluation (immediate)	Lazy evaluation (deferred)
Parallel Processing	Limited	Supports multi-threading
Integration	Seamless with other Python libraries	Integrates well with common data science libraries
API	Intuitive and widely adopted	Similar to Pandas, easy to learn
Feature Set	Extensive (data manipulation, analysis, etc.)	Comparable for large datasets, includes virtual columns, efficient filtering, and binning
Memory Usage	High for large datasets	Low due to memory-mapped files
Installation	<code>pip install pandas</code>	<code>pip install vaex</code>
Documentation	Extensive and well-documented	Extensive and well-documented
Community Support	Large and active community	Growing and active community

Source 1: <https://www.datacamp.com/tutorial/benchmarking-high-performance-pandas-alternatives#>

Source 2: <https://medium.com/helpshift-engineering/vaex-pandas-on-steroids-ded8ca765f55>

Outline



Setting the Context



Introduction to advanced Python libraries



High-performance data manipulation with Vaex



Interactive data exploration with D-Tale



Feature engineering with Feature Engine

Importance of Advanced Python Libraries in Enhancing Data Analysis Capabilities and Improving Workflow Efficiency

1. Efficiency and Performance:

•Handling Large Datasets:

- Advanced libraries like Vaex are designed to work with large datasets that do not fit into memory, enabling efficient out-of-core data manipulation and analysis. This capability is crucial for processing large-scale data in a timely manner.

•Optimized Computations:

- These libraries leverage optimized algorithms and data structures, significantly reducing computation time. For example, Vaex uses memory-mapped files and efficient indexing to speed up data operations.

•Parallel Processing:

- Many advanced libraries support parallel processing and distributed computing, allowing for faster data processing by utilizing multiple CPU cores or even clusters of machines.

2. Ease of Use and Flexibility:

•User-Friendly APIs:

- Libraries like D-Tale provide intuitive interfaces for data exploration and

visualization, making it easier for data professionals to interact with their data without extensive coding.

•Seamless Integration:

- These libraries integrate seamlessly with other popular Python libraries such as Pandas, NumPy, and scikit-learn, allowing for a smooth transition and interoperability between different tools in the data science ecosystem.

3. Community Support and Documentation:

•Comprehensive Documentation:

- Most advanced libraries come with extensive documentation and tutorials, providing clear guidelines on how to use their features effectively. This makes it easier for users to get started and overcome any challenges they may encounter.

•Active Community:

- An active community of users and developers contributes to continuous improvement and support. Users can seek help from community forums, GitHub issues, and other online resources, ensuring they are not alone in their learning journey.

4. Improving Workflow Efficiency:

•Streamlined Data Processing:

- Feature Engine simplifies common preprocessing tasks such as handling missing data, encoding categorical variables, and scaling features, making the data preparation phase more efficient and less error-prone.

•Interactive Analysis:

- D-Tale provides an interactive environment for data exploration, enabling quick insights and facilitating a more iterative and exploratory analysis process. This leads to faster discovery of patterns and anomalies in the data.

5. Enhancing Analytical Capabilities:

•Advanced Techniques:

- By leveraging advanced libraries, data professionals can implement sophisticated techniques and models that were previously challenging or time-consuming. This includes out-of-core algorithms, automated machine learning, and complex feature engineering.

•Better Decision Making:

- The ability to analyze large datasets efficiently and accurately leads to more informed and data-driven decisions. This is particularly important in industries like healthcare, where timely and accurate insights can significantly impact business outcomes.

•Scalability:

- Advanced libraries are designed to scale with the data. As the volume and

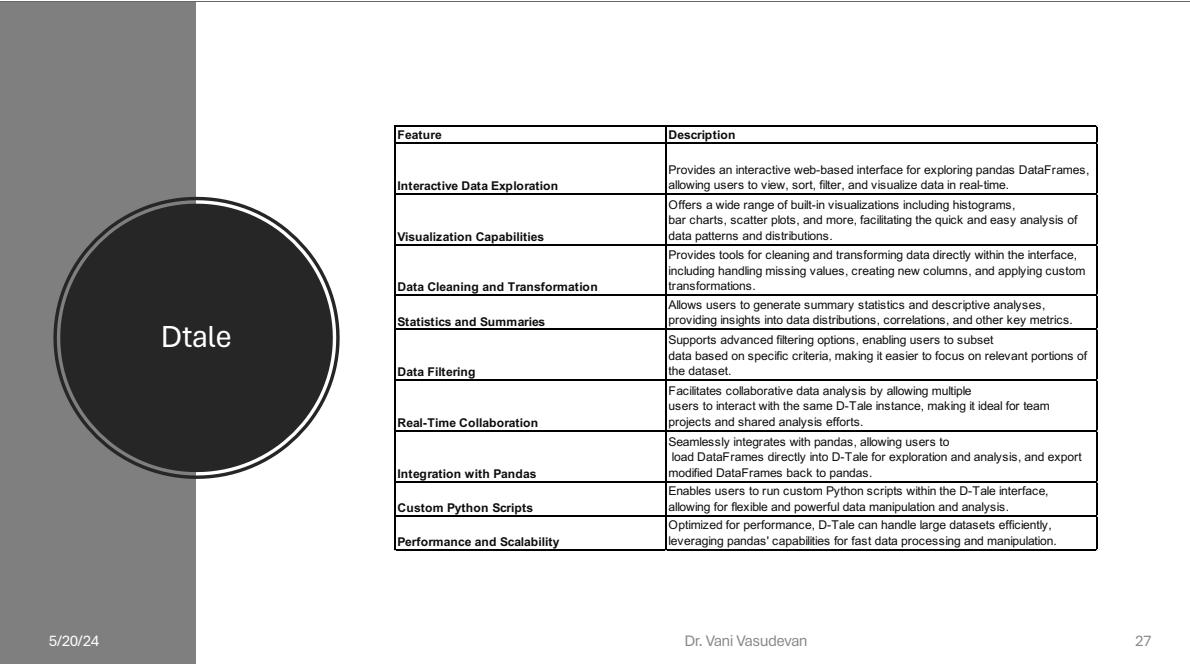
complexity of data grow, these tools ensure that data analysis remains manageable and performant.

Dtale

D-Tale is a Python library that provides an interactive and web-based interface for exploring and visualizing datasets.

It is the combination of a Flask back-end and a React front-end to bring you an easy way to view & analyze Panda's data structures.

26



Outline



Setting the Context



Introduction to advanced Python libraries



High-performance data manipulation with Vaex



Interactive data exploration with D-Tale



Feature engineering with Feature Engine

Importance of Advanced Python Libraries in Enhancing Data Analysis Capabilities and Improving Workflow Efficiency

1. Efficiency and Performance:

•Handling Large Datasets:

- Advanced libraries like Vaex are designed to work with large datasets that do not fit into memory, enabling efficient out-of-core data manipulation and analysis. This capability is crucial for processing large-scale data in a timely manner.

•Optimized Computations:

- These libraries leverage optimized algorithms and data structures, significantly reducing computation time. For example, Vaex uses memory-mapped files and efficient indexing to speed up data operations.

•Parallel Processing:

- Many advanced libraries support parallel processing and distributed computing, allowing for faster data processing by utilizing multiple CPU cores or even clusters of machines.

2. Ease of Use and Flexibility:

•User-Friendly APIs:

- Libraries like D-Tale provide intuitive interfaces for data exploration and

visualization, making it easier for data professionals to interact with their data without extensive coding.

•Seamless Integration:

- These libraries integrate seamlessly with other popular Python libraries such as Pandas, NumPy, and scikit-learn, allowing for a smooth transition and interoperability between different tools in the data science ecosystem.

3. Community Support and Documentation:

•Comprehensive Documentation:

- Most advanced libraries come with extensive documentation and tutorials, providing clear guidelines on how to use their features effectively. This makes it easier for users to get started and overcome any challenges they may encounter.

•Active Community:

- An active community of users and developers contributes to continuous improvement and support. Users can seek help from community forums, GitHub issues, and other online resources, ensuring they are not alone in their learning journey.

4. Improving Workflow Efficiency:

•Streamlined Data Processing:

- Feature Engine simplifies common preprocessing tasks such as handling missing data, encoding categorical variables, and scaling features, making the data preparation phase more efficient and less error-prone.

•Interactive Analysis:

- D-Tale provides an interactive environment for data exploration, enabling quick insights and facilitating a more iterative and exploratory analysis process. This leads to faster discovery of patterns and anomalies in the data.

5. Enhancing Analytical Capabilities:

•Advanced Techniques:

- By leveraging advanced libraries, data professionals can implement sophisticated techniques and models that were previously challenging or time-consuming. This includes out-of-core algorithms, automated machine learning, and complex feature engineering.

•Better Decision Making:

- The ability to analyze large datasets efficiently and accurately leads to more informed and data-driven decisions. This is particularly important in industries like healthcare, where timely and accurate insights can significantly impact business outcomes.

•Scalability:

- Advanced libraries are designed to scale with the data. As the volume and

complexity of data grow, these tools ensure that data analysis remains manageable and performant.

Feature-Engine

Feature-engine is a Python library with multiple transformers to engineer and select features for machine learning models. Feature-engine adopts Scikit-learn functionality with methods `fit()` and `transform()` to learn parameters from and then transform the data.

29

Feature-engine

Feature	Description
Transformations for Feature Engineering	Provides a variety of transformers for feature engineering, including handling missing values, encoding categorical variables, and transforming numerical variables.
Pipeline Integration	Designed to integrate seamlessly with Scikit-Learn pipelines, allowing for easy incorporation of feature engineering steps in machine learning workflows.
Handling Missing Data	Includes transformers for imputing missing data using mean, median, mode, arbitrary values, or methods based on data characteristics.
Encoding Categorical Variables	Offers various encoding methods such as one-hot encoding, ordinal encoding, count or frequency encoding, target mean encoding, and more.
Variable Transformation	Allows for transformation of numerical variables through logarithmic, reciprocal, power, and box-cox transformations, among others.
Discretization	Provides methods for discretizing continuous variables into discrete intervals or bins, which can be useful for certain types of models and analyses.
Feature Scaling	Implements scaling techniques such as standardization, min-max scaling, and robust scaling to normalize data.
Feature Selection	Includes techniques for selecting relevant features based on statistical tests, model importance, or user-defined criteria.
Combining Variables	Allows for the creation of new features by combining existing variables through mathematical operations or aggregations.
Outlier Handling	Provides methods to cap, transform, or remove outliers from the data, ensuring that extreme values do not skew the model performance.

Thank you

Any Questions?

Reach me @ vani.v@nmit.ac.in

LinkedIn: www.linkedin.com/in/dr-vani-vasudevan-0b89b713

GitHub:
<https://github.com/vanivasudevan/AdvancedLibrariesDemo>

