

Introduction to Basic Data Analytics Tools

What is Data Analytics?

Data analytics is the science of analyzing raw data in order to make conclusions about that information.

Data Analytics Pipeline



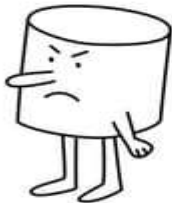
Data Acquisition

How To Collect Data!

Unstructured Data



Traditional RDBMS



© 2011 TimeEdScott.com

"I'm sorry, I'm just not that into you..."

- ★ REST API
- ★ From end users
- ★ Web scrape
- ★ Email and cloud storage
- ★ Client's server

REST : Representational State Transfer

POST | PUT | Get

Requests

A library for making HTTP requests in Python.

Key Features:

- Keep-alive & Connection Pooling
- Sessions with Cookie Persistence

```
import requests

url = 'https://jsonplaceholder.typicode.com/posts/1'
response = requests.get(url=url)
print(response.text)

payload = {
    "userID": "1",
    "title": "thisTitle",
    "body": "thisBody"
}

url = 'https://jsonplaceholder.typicode.com/posts'
response = requests.post(url=url, data=payload)

print(response.status_code)
print(response.text)
```

BeautifulSoup

A library for parsing HTML and XML documents.

Key Features:

- Multiple parser support (e.g., lxml, html5lib, and others)
- Creates parse tree which is easy to navigate

```
from bs4 import BeautifulSoup

html_doc = """
<html>
  <head>
    <title>Example HTML Document</title>
  </head>
  <body>
    <p class="title">Example paragraph</p>
    <a href="https://jsonplaceholder.typicode.com/posts/1">Link 1</a>
    <a href="https://jsonplaceholder.typicode.com/posts/2">Link 2</a>
  </body>
</html>
"""

soup = BeautifulSoup(html_doc, 'lxml')

print(soup.title)
# Output: <title>Example HTML Document</title>

print(soup.find_all('a')[0]['href'])
# Output: https://jsonplaceholder.typicode.com/posts/1
```

Flask, Flask-RESTPlus and Swagger UI

Flask is a micro web framework written in Python.

Flask-RESTPlus is an extension for Flask that adds support for quickly building REST APIs. It automatically documents the APIs which is visible in Swagger UI.



The screenshot shows the Swagger UI interface for an API. At the top left, it displays "API" with a version number "1.0.0" and the base URL "http://localhost:5001/swagger-ui". A green "Authorize" button is located in the top right corner. Below this, there is a section for "Posts" APIs, which is expanded to show two endpoints: a GET endpoint for "/posts/" and a POST endpoint for "/posts/". Each endpoint has a lock icon indicating it is protected. At the bottom, there is a "Models" section with a right-pointing arrow.

Data Pre-Processing and Storage

How To Clean Data!

- ★ Remove duplicate
- ★ Validate
- ★ Handle missing data
- ★ Fix errors
- ★ Filter outliers

How To Store Data!

- ★ RDBMS
- ★ ORM



*"What's a data lake for?
So you can drown in more data even faster!"*

Pandas

A library for data manipulation and analysis.

Key Features:

- Loading data into in-memory data objects from different file formats.
- Data alignment and integrated handling of missing data.

```
import pandas as pd
from fuzzywuzzy import process

df = pd.DataFrame(
    columns=['Division', 'Code'],
    data=[['Dhak', '30'], ['Chattogram', '29']]
)

divisions = ['Dhaka', 'Chattogram']

for i in range(len(df)):
    division = process.extractOne(df.Division[i], divisions)[0]
    df.at[i, 'Division'] = division

print(df)
#      Division Code
# 0      Dhaka    30
# 1  Chattogram    29
```

SQLAlchemy

SQLAlchemy is a popular SQL toolkit and Object Relational Mapper.

Key Features:

- Function-based query construction.
- Multiple database support (e.g., SQLite, Postgresql, MySQL, Oracle, MS-SQL, Firebird, Sybase and others).

```
from sqlalchemy import create_engine, MetaData, Table, Column, Integer, String
from sqlalchemy.sql.functions import user
import os

path = os.path.dirname(__file__)
engine = create_engine('sqlite:///{}'.format(os.path.join(path, 'database.db')))
meta = MetaData()

posts = Table(
    'posts', meta,
    Column('id', Integer, primary_key=True, autoincrement=True),
    Column('userID', Integer),
    Column('title', String),
    Column('body', String),
)

meta.create_all(engine)

conn = engine.connect()

result = conn.execute(posts.insert(), [
    {'userID': 0, 'title': 'title 0', 'body': 'body 0'},
    {'userID': 1, 'title': 'title 1', 'body': 'body 1'},
])

result = conn.execute(posts.select())
for row in result:
    print(row)
```

Data Analysis

How To Analyze Data!



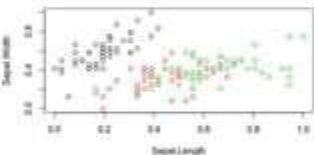
"If you don't reveal some insights soon, I'm going to be forced to slice, dice, and drill!"

- ★ Five number summary
(maximum, minimum, median, 1st quartile, 3rd quartile)
- ★ Average
- ★ Standard Deviation
- ★ Ratio
- ★ Interval
- ★ Trends
- ★ Aggregate and group by
- ★ Regression
- ★ Clustering

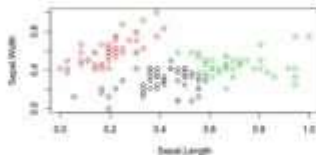
R and RStudio

R is a popular programming language for data analysis. RStudio is an IDE for R.

Original Classes



Clusters by k-means



```
require("datasets")
data("iris")
head(iris)
summary(iris)

X <- iris[,c(1, 2)]
y <- iris[, "Species"]

normalize <- function(x) {
  return ((x - min(x)) / (max(x) - min(x)))
}

X$Sepal.Length <- normalize(X$Sepal.Length)
X$Sepal.Width <- normalize(X$Sepal.Width)
head(X)

result <- kmeans(X, 3)

result$size
result$centers

plot(X, col=result$cluster)
plot(X, col=y)
```

Data Presentation

How To Visualize Data!

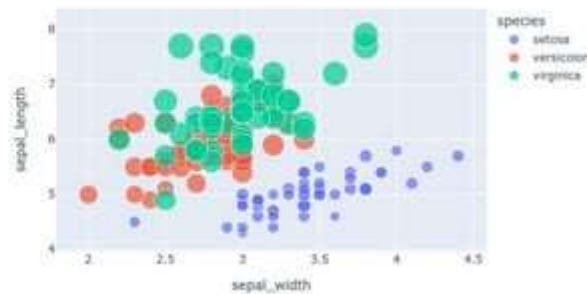


"Fake news!"

- ★ Charts
 - Line
 - Bar
 - Pie
 - Scatter
- ★ Graphs
- ★ Maps
 - Bubble
 - Polygon
- ★ Dashboards

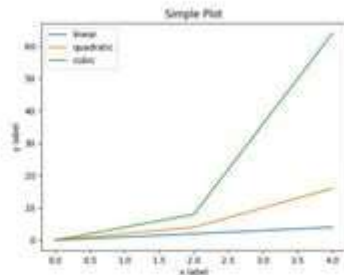
Plotly

An interactive graphing library.



Matplotlib

A plotting library for Python.

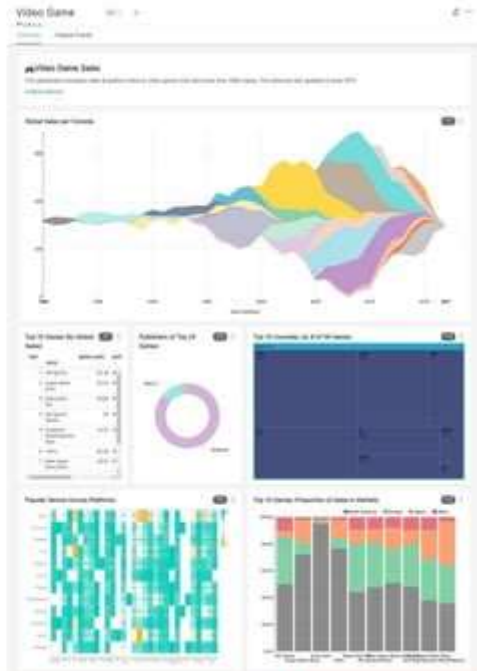


Apache Superset

A Data Visualization and Data Exploration Platform.

Key Features:

- It supports all the data sources that support SQL Alchemy and supports querying using SQL.
- Superset allows sharing dashboards.
- It comes with security features like Authentication, User Management and Roles.



Other Notable Tools

- ★ Excel
- ★ Tableau Public
- ★ Grafana
- ★ Microsoft Power BI
- ★ And many more . . .

Challenges

1. Poor quality data
2. Data privacy and security
3. Weak infrastructure
4. Data from multiple sources
5. Scaling data analysis



"Sir, some citizens have been complaining about our data privacy policies..."

*Do you want their names?
Social security numbers?
Most embarrassing secrets?"*

Links

Flask: <https://flask.palletsprojects.com/en/2.0.x/>

Flask-RESTPlus: <https://flask-restplus.readthedocs.io/en/stable/>

SQLAlchemy: <https://www.sqlalchemy.org/>

R Programming Language: <https://www.r-project.org/>

k-means Clustering: https://en.wikipedia.org/wiki/K-means_clustering

Plotly: <https://plotly.com/>

Superset Docs: <https://superset.apache.org/docs/intro>

Presentation GitHub Link: <https://github.com/saadrumon/basic-data-analytics-tools-presentation.git>

Thank You

Any Questions?

“Information is the oil of the 21st century, and analytics is the combustion engine.”

- Peter Sondergaard