# 資料探勘
# (Data Mining)

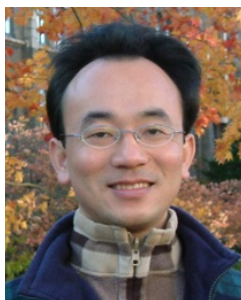# 資料科學與資料探勘：
# 發現，分析，可視化和呈現數據

**(Data Science and Data Mining: Discovering, Analyzing, Visualizing and Presenting Data)**

**Min-Yuh Day**
**戴敏育**
**Associate Professor**
副教授
**Institute of Information Management**, **National Taipei University**
**國立臺北大學 資訊管理研究所**

https://web.ntpu.edu.tw/~myday

2021-03-16

# 課程大綱 (Syllabus)

週次 (Week)　日期 (Date)　內容 (Subject/Topics)

1  2021/02/23  資料探勘介紹 (Introduction to data mining)

2  2021/03/02  ABC：人工智慧，大數據，雲端運算
　　　　　　　(ABC: AI, Big Data, Cloud Computing)

3  2021/03/09  Python資料探勘的基礎
　　　　　　　(Foundations of Data Mining in Python)

4  2021/03/16  資料科學與資料探勘：發現，分析，可視化和呈現數據
　　　　　　　(Data Science and Data Mining:
　　　　　　　 Discovering, Analyzing, Visualizing and Presenting Data)

5  2021/03/23  非監督學習：關聯分析，購物籃分析
　　　　　　　 (Unsupervised Learning: Association Analysis,
　　　　　　　　Market Basket Analysis)

6  2021/03/30  資料探勘個案研究 I
　　　　　　　 (Case Study on Data Mining I)

# 課程大綱 (Syllabus)

週次 (Week)　日期 (Date)　內容 (Subject/Topics)

7  2021/04/06 非監督學習：集群分析，行銷市場區隔
(Unsupervised Learning: Cluster Analysis, Market Segmentation)

8  2021/04/13 監督學習：分類和預測
(Supervised Learning: Classification and Prediction)

9  2021/04/20 期中報告 (Midterm Project Report)

10  2021/04/27 監督學習：分類和預測
(Supervised Learning: Classification and Prediction)

11  2021/05/04 機器學習和深度學習
(Machine Learning and Deep Learning)

12  2021/05/11 卷積神經網絡
(Convolutional Neural Networks)

# 課程大綱 (Syllabus)

週次 (Week)　日期 (Date)　內容 (Subject/Topics)

13　2021/05/18　資料探勘個案研究 II
　　　　　　　　(Case Study on Data Mining II)

14　2021/05/25　遞歸神經網絡
　　　　　　　　(Recurrent Neural Networks)

15　2021/06/01　強化學習
　　　　　　　　(Reinforcement Learning)

16　2021/06/08　社交網絡分析
　　　　　　　　(Social Network Analysis)

17　2021/06/15　期末報告 I (Final Project Report I)

18　2021/06/22　期末報告 II (Final Project Report II)
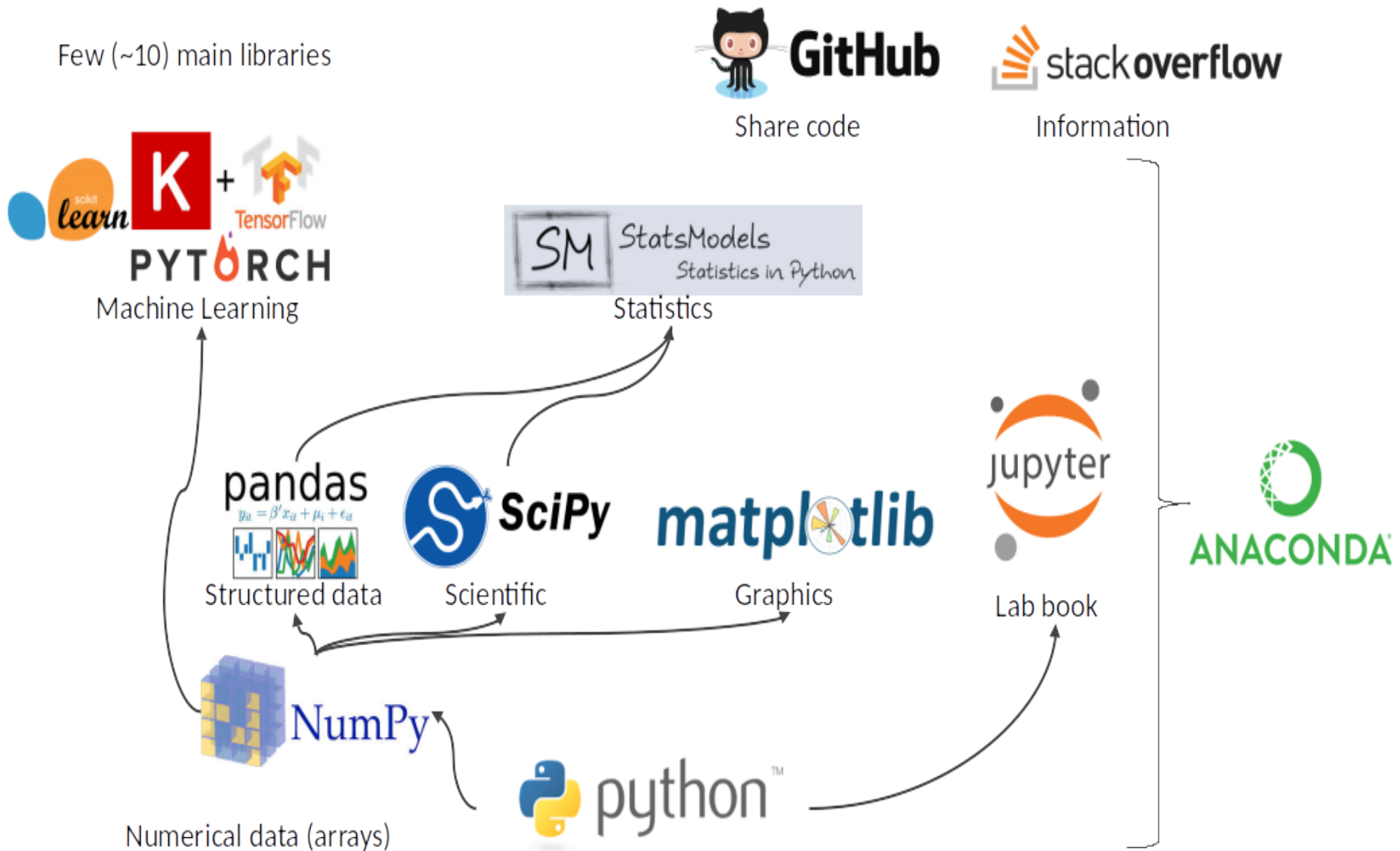
# Data Science and Data Mining: Discovering, Analyzing, Visualizing and Presenting Data
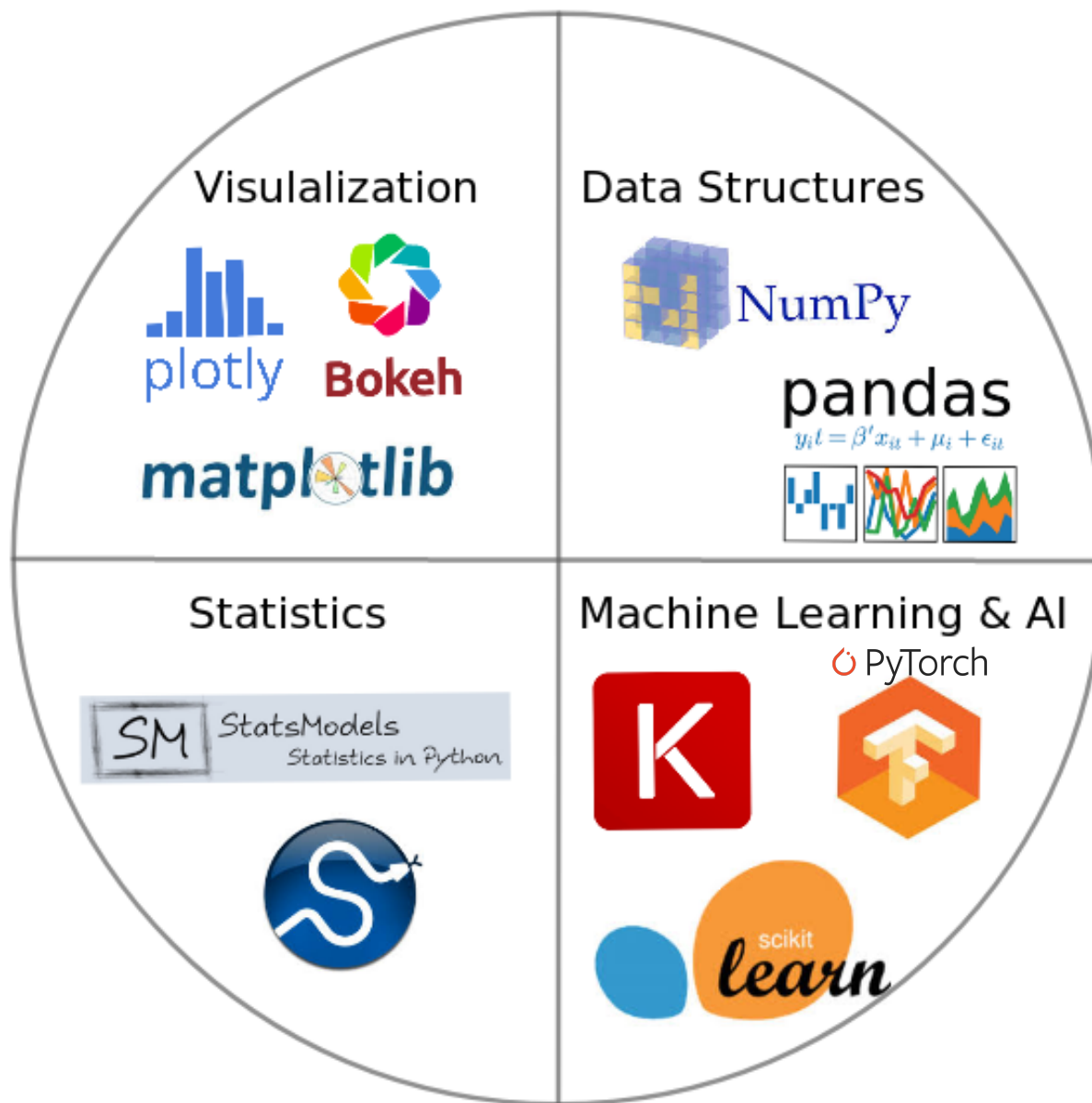
# Outline

- **Data Science and Data Mining**

- **Discovering, Analyzing, Visualizing and Presenting Data with Python**
  - **Pandas**

  - **Matplotlib**

  - **Seaborn**
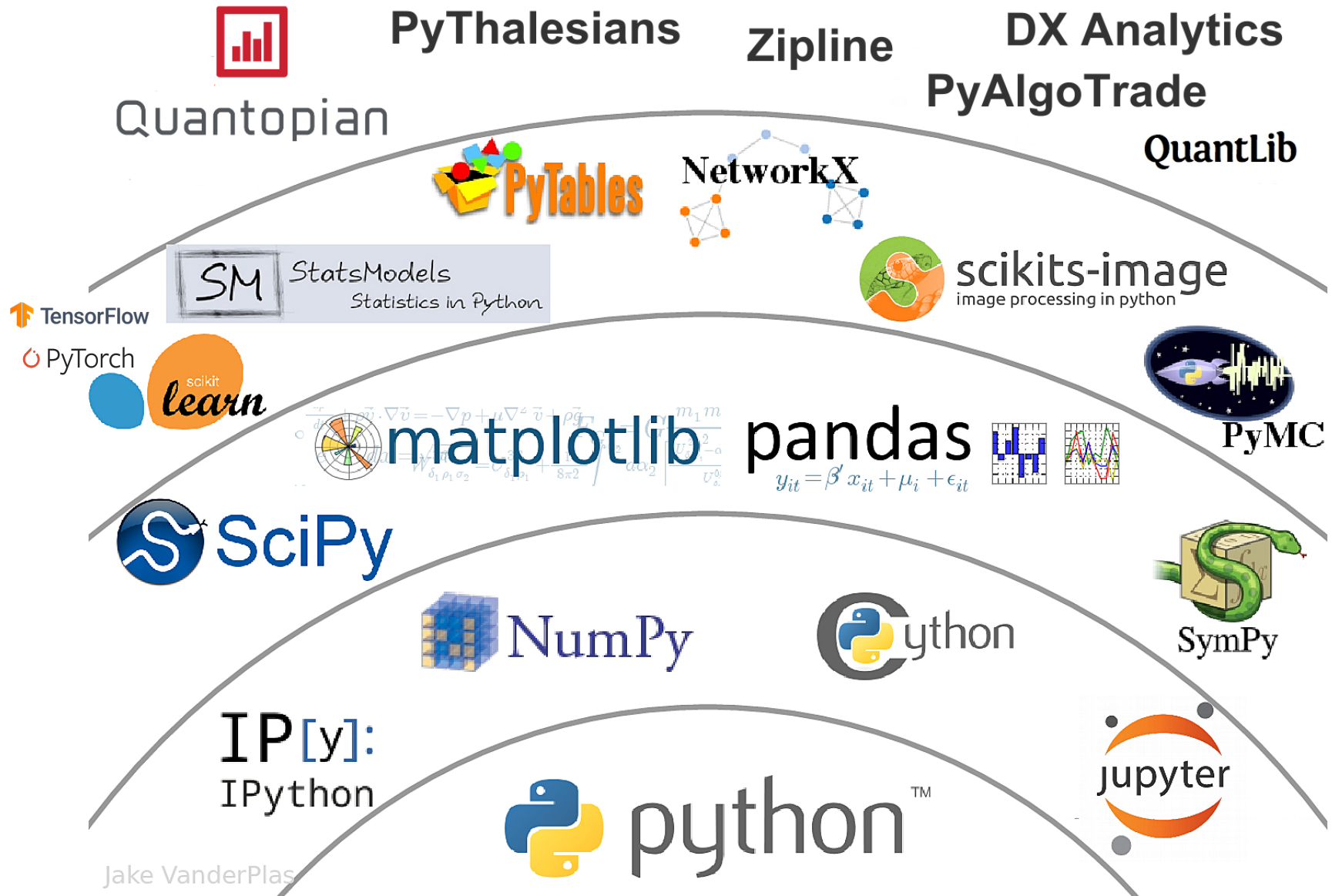
  - **Plotly**

  - **Bokeh, Altair**
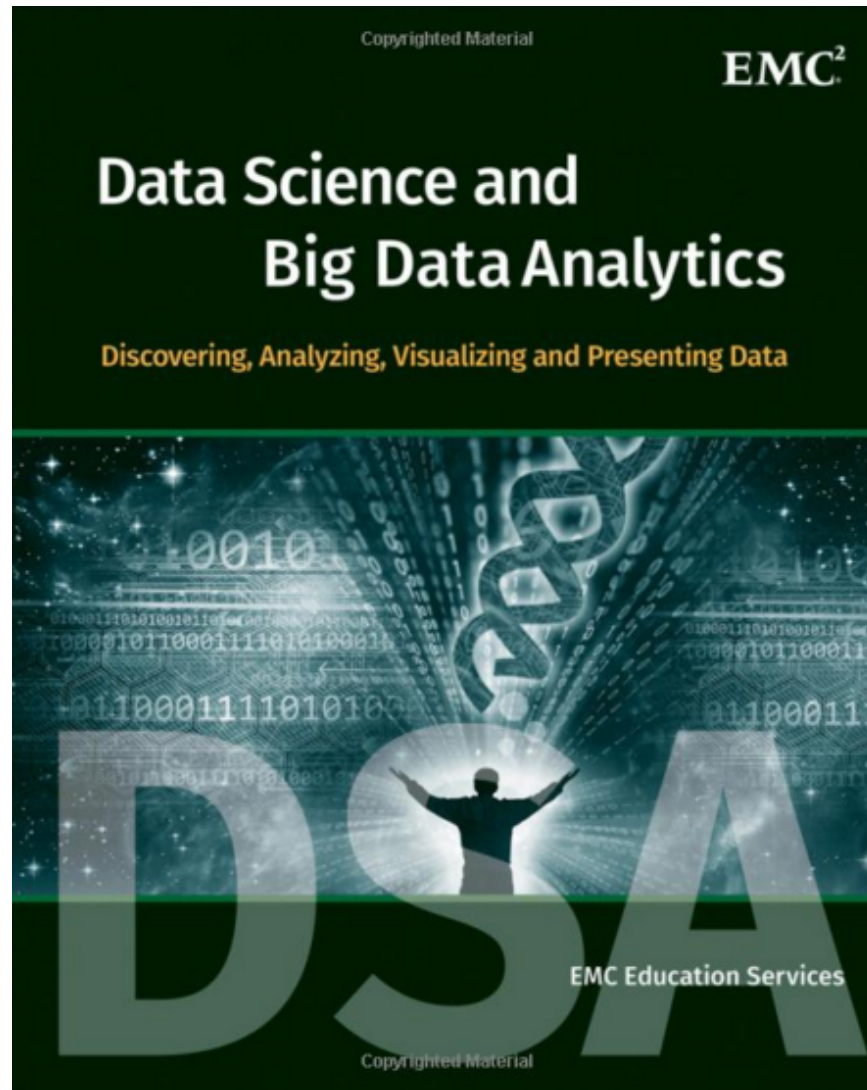
# Python Ecosystem for Data Science

Few (~10) main libraries

GitHub
Share code

stackoverflow
Information

scikit learn · K + TensorFlow · PYTORCH
Machine Learning

SM StatsModels Statistics in Python
Statistics

pandas $y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$
Structured data

SciPy
Scientific

matplotlib
Graphics

jupyter
Lab book

ANACONDA

NumPy
Numerical data (arrays)

python

# Python Ecosystem for Data Science

# The Quant Finance PyData Stack

# EMC Education Services,
## Data Science and Big Data Analytics:
## Discovering, Analyzing, Visualizing and Presenting Data,
## Wiley, 2015

# Data Science

# Data Analyst

- Data analyst is just another term for professionals who were doing BI in the form of data compilation, cleaning, reporting, and perhaps some visualization.

- Their skill sets included Excel, some SQL knowledge, and reporting.

- You would recognize those capabilities as descriptive or reporting analytics.

# Data Scientist

- Data scientist is responsible for predictive analysis, statistical analysis, and more advanced analytical tools and algorithms.

- They may have a deeper knowledge of algorithms and may recognize them under various labels—data mining, knowledge discovery, or machine learning.

- Some of these professionals may also need deeper programming knowledge to be able to write code for data cleaning/analysis in current Web-oriented languages such as Java or Python and statistical languages such as R.

- Many analytics professionals also need to build significant expertise in statistical modeling, experimentation, and analysis.

# Data Science and Business Intelligence



Predictive Analytics and Data Mining (Data Science)

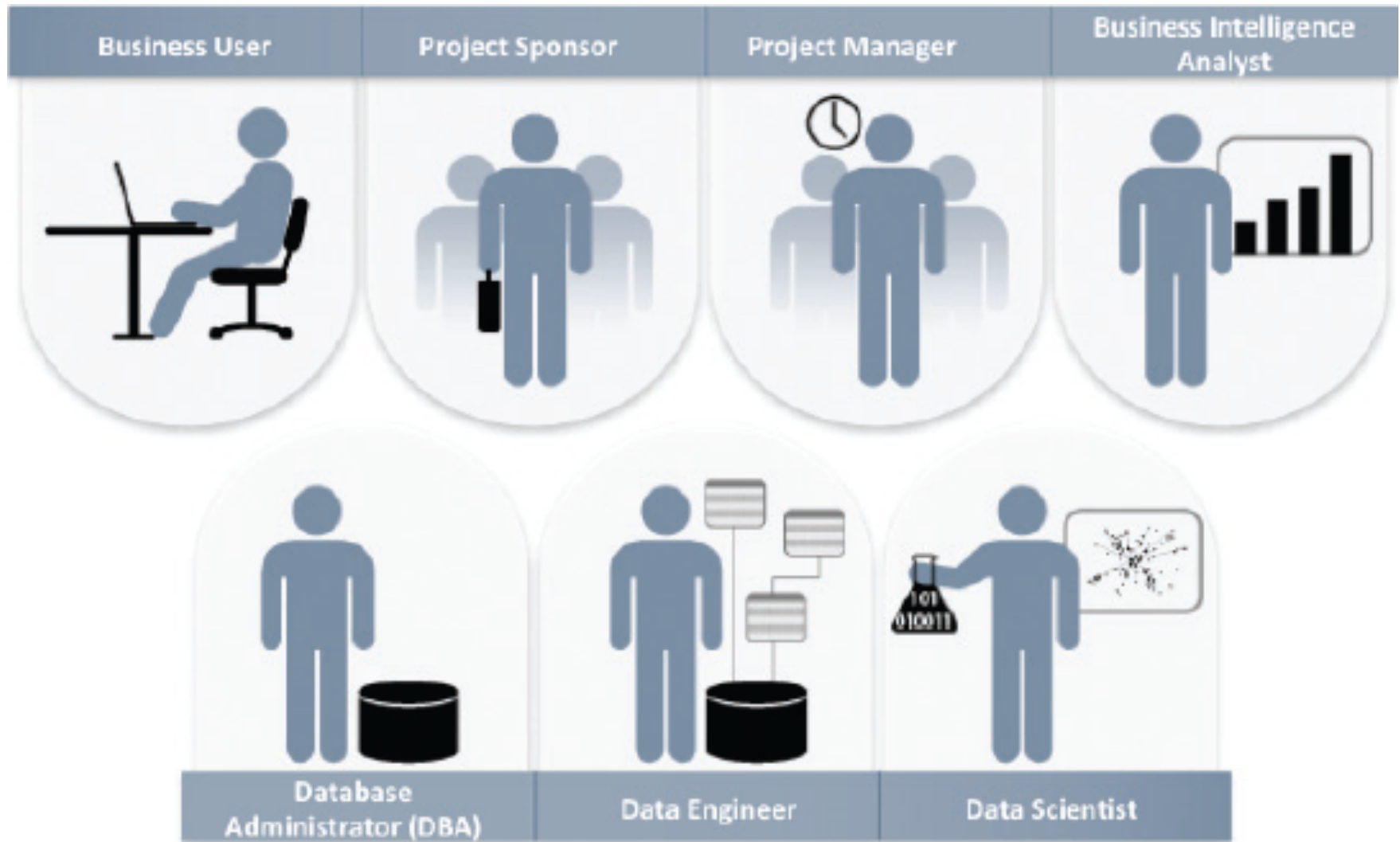| Typical Techniques and Data Types | • Optimization, predictive modeling, forecasting, statistical analysis<br>• Structured/unstructured data, many types of sources, very large datasets |
|---|---|
| Common Questions | • What if...?<br>• What's the optimal scenario for our business?<br>• What will happen next? What if these trends continue? Why is this happening? |

Business Intelligence

| Typical Techniques and Data Types | • Standard and ad hoc reporting, dashboards, alerts, queries, details on demand<br>• Structured data, traditional sources, manageable datasets |
|---|---|
| Common Questions | • What happened last quarter?<br>• How many units sold?<br>• Where is the problem? In which situations? |

# Data Science and Business Intelligence



Exploratory

**Predictive Analytics and Data Mining (Data Science)**

| Typical Techniques and Data Types | • Optimization, predictive modeling, forecasting, statistical analysis<br>• Structured/unstructured data, many types of sources, very large datasets |
|---|---|
| Common Questions | • What if...?<br>• What's the optimal scenario for our business?<br>• What will happen next? What if these trends continue? Why is this happening? |

# Predictive Analytics and Data Mining (Data Science)

Past

Time

Future

# Predictive Analytics and Data Mining
## (Data Science)

Structured/unstructured data, many types of sources, very large datasets

Optimization, predictive modeling, forecasting statistical analysis

What if…?
What's the optimal scenario for our business?
What will happen next?
What if these trends countinue?
Why is this happening?

# Profile of a Data Scientist

- **Quantitative**
  - mathematics or statistics
- **Technical**
  - software engineering, machine learning, and programming skills
- **Skeptical mind-set** and **critical thinking**
- **Curious** and **creative**
- **Communicative** and **collaborative**

# Data Scientist Profile

# Big Data Analytics Lifecycle

# Key Roles for a Successful Analytics Project

# Overview of Data Analytics Lifecycle

# Overview of Data Analytics Lifecycle

1. Discovery

2. Data preparation

3. Model planning

4. Model building

5. Communicate results

6. Operationalize

# Key Outputs from a Successful Analytics Project



Source: EMC Education Services, Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data, Wiley, 2015

# Example of Analytics Applications in a Retail Value Chain

## Retail Value Chain
Critical needs at every touch point of the Retail Value Chain



- Shelf-space optimization
- Location analysis
- Shelf and floor planning
- Promotions and markdown optimization

- Trend analysis
- Category management
- Predicting trigger events for sales
- Better forecasts of demand

- Deliver seamless customer experience
- Understand relative performance of channels
- Optimize marketing strategies

**Vendors** → Planning → Merchandizing → Buying → Warehouse & Logistics → Multichannel Operations → **Customers**

- Supply chain management
- Inventory cost optimization
- Inventory shortage and excess management
- Less unwanted costs

- Targeted promotions
- Customized inventory
- Promotions and price optimization
- Customized shopping experience

- On-time product availability at low costs
- Order fulfillment and clubbing
- Reduced transportation costs

- Building retention and satisfaction
- Understanding the needs of the customer better
- Serving high LTV customers better

# Analytics Ecosystem

# Job Titles of Analytics

Source: Ramesh Sharda, Dursun Delen, and Efraim Turban (2017), Business Intelligence, Analytics, and Data Science: A Managerial Perspective, 4th Edition, Pearson

# Three Types of Analytics

**Business Analytics**

| Descriptive | Predictive | Prescriptive |
|---|---|---|

| | Descriptive | Predictive | Prescriptive |
|---|---|---|---|
| **Questions** | What happened? What is happening? | What will happen? Why will it happen? | What should I do? Why should I do it? |
| **Enablers** | ✓ Business reporting<br>✓ Dashboards<br>✓ Scorecards<br>✓ Data warehousing | ✓ Data mining<br>✓ Text mining<br>✓ Web/media mining<br>✓ Forecasting | ✓ Optimization<br>✓ Simulation<br>✓ Decision modeling<br>✓ Expert systems |
| **Outcomes** | **Well-defined business problems and opportunities** | **Accurate projections of future events and outcomes** | **Best possible business decisions and actions** |

# A Data to Knowledge Continuum

# A Simple Taxonomy of Data

# Data Preprocessing Steps

# An Analytics Approach to Predicting Student Attrition

# A Graphical Depiction of the Class Imbalance Problem

# Relationship between Statistics and Descriptive Analytics

# Understanding the Specifics about Box-and-Whiskers Plots

# Relationship between Dispersion and Shape Properties.

# A Scatter Plot and a Linear Regression Line

# A Process Flow for Developing Regression Models.

# The Logistic Function

$$f(y) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

# Predicting NCAA Bowl Game Outcomes

Source: Ramesh Sharda, Dursun Delen, and Efraim Turban (2017), Business Intelligence, Analytics, and Data Science: A Managerial Perspective, 4th Edition, Pearson

# A Sample Time Series of Data on Quarterly Sales Volumes



Quarterly Product Sales (in Millions)

# The Role of Information Reporting in Managerial Decision Making

# A Taxonomy of Charts and Graphs



Two variables per item
Many categories
Many items
Few items
Few categories
One variable per item
Among items
Cyclic data
Noncyclic data
Many periods
Over time
Single or few categories
Many categories
Few periods

Comparison

Two variables
Three variables
Relationship

**What would you like to show in your chart or graph?**

Distribution
Single variable
Few data points
Many data point
Two variables
Three variables

Composition

Changing over time
Static

Few periods
Many periods

Only relative difference matters
Relative and absolute difference matters
Only relative difference matters
Relative and absolute difference matters
Simple share of total
Accumulation or subtraction to total
Components of components

# A Gapminder Chart That Shows the Wealth and Health of Nations

# Magic Quadrant for Business Intelligence and Analytics Platforms



Source: https://www.tableau.com/reports/gartner

44

# A Storyline Visualization in Tableau Software

# An Overview of SAS Visual Analytics Architecture

# A Screenshot from SAS Visual Analytics

# A Sample Executive Dashboard

# Exploratory Network Analysis



**1** see the network

1st graph viz tool: Pajek (1996)
Vladimir Batagelj, Andrej Mrvar

**2** interact in real time

Gephi prototype (2008)
group, filter, compute metrics...

**3** build a visual language

size by rank, color by partition,
label, curved edges, thickness...

A—B
B—A

Source: http://sebastien.pro/gephi-icwsm-tutorial.pdf

# Looking for a "Simple Small Truth"? What Data Visualization Should Do?



1. Make complex things **simple**
2. Extract **small** information from large data
3. Present **truth**, do not deceive

Source: http://sebastien.pro/gephi-icwsm-tutorial.pdf

# Gephi

# Discovering, Analyzing, Visualizing and Presenting Data with Python in Google Colab

# Python Data Analysis and Visualization

# Python

# Pandas

# Python
# matplotlib

# Python
# seaborn

# Python

# bokeh

# Python
# Altair

# Python matplotlib

# Python Seaborn

# seaborn: statistical data visualization

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

For a brief introduction to the ideas behind the library, you can read the introductory notes. Visit the installation page to see how you can download the package and get started with it. You can browse the example gallery to see what you can do with seaborn, and then check out the tutorial and API reference to find out how.

To see the code or report a bug, please visit the GitHub repository. General support questions are most at home on stackoverflow or discourse, which have dedicated channels for seaborn.

## Contents

- Introduction
- Release notes
- Installing
- Example gallery
- Tutorial
- API reference

## Features

- Relational: API | Tutorial
- Distribution: API | Tutorial
- Categorical: API | Tutorial
- Regression: API | Tutorial
- Multiples: API | Tutorial
- Style: API | Tutorial
- Color: API | Tutorial

Back to top

https://seaborn.pydata.org/

61

# Python Plotly Graphing Library

**plotly | Graphing Libraries**

GitHub Star 9,085    **DO MORE WITH DASH**

**Quick Start**

Getting Started

Is Plotly Free?

Figure Reference

API Reference

Dash

GitHub

community.plotly.com

**Examples**

Fundamentals

Basic Charts

Statistical Charts

Artificial Intelligence and Machine Learning

Scientific Charts

Financial Charts

Maps
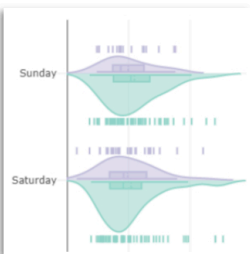
3D Charts

## Plotly Python Open Source Graphing Library

Plotly's Python graphing library makes interactive, publication-quality graphs. Examples of how to make line plots, scatter plots, area charts, bar charts, error bars, box plots, histograms, heatmaps, subplots, multiple-axes, polar charts, and bubble charts.
Plotly.py is free and open source and you can view the source, report issues or contribute on GitHub.

> *Our recommended IDE for Plotly's Python graphing library is Dash Enterprise's Data Science Workspaces, which has both Jupyter notebook and Python code file support.*
> *Find out if your company is using Dash Enterprise.*
>
> **Install Dash Enterprise on Azure** | **Install Dash Enterprise on AWS**

## Fundamentals

More Fundamentals »

The Figure Data Structure

Creating and Updating Figures

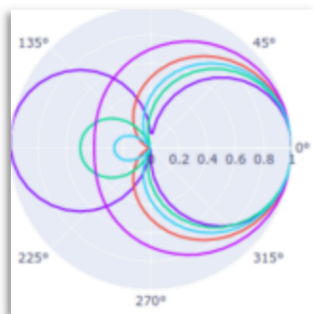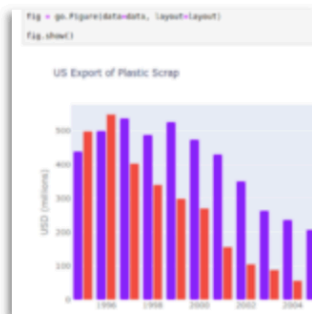Displaying Figures

Plotly Express

Analytical Apps with Dash

https://plotly.com/python/

**plotly**

# Python Plotly Graphing Library

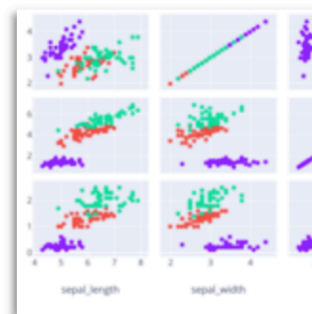## Fundamentals

More Fundamentals »



The Figure Data Structure



Creating and Updating Figures



Displaying Figures


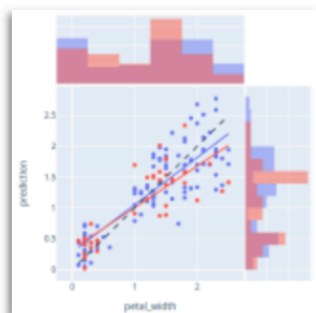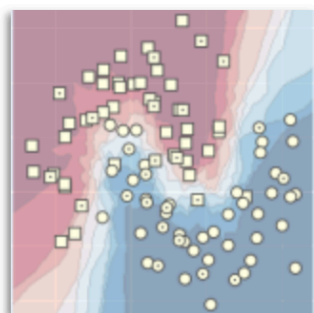
Plotly Express



Analytical Apps with Dash

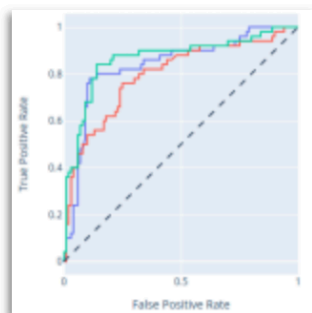## Artificial Intelligence and Machine Learning

More AI and ML »



ML Regression



kNN Classification
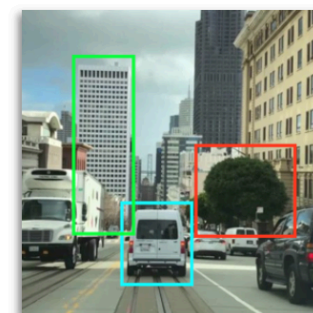


ROC and PR Curves



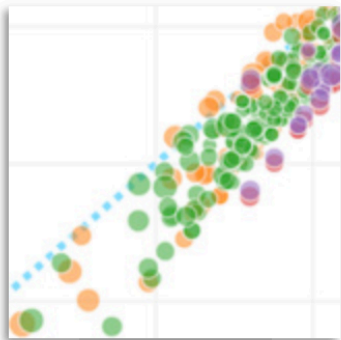PCA Visualization



AI/ML Apps with Dash

https://plotly.com/python/
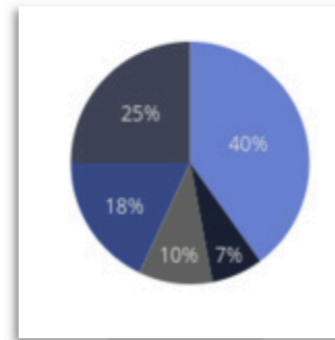
# Python Plotly Graphing Library

## Basic Charts

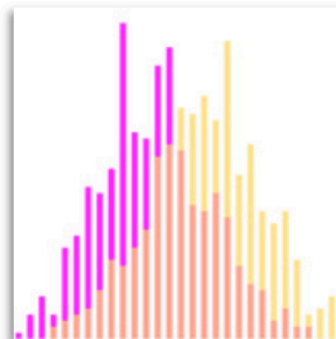| Scatter Plots | Line Charts | Bar Charts | Pie Charts | Bubble Charts |

## Statistical Charts

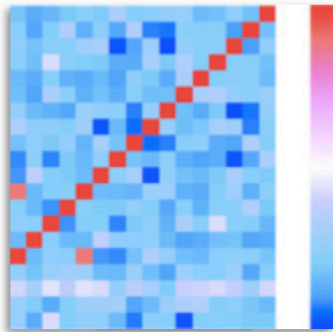| Error Bars | Box Plots | Histograms | Distplots | 2D Histograms |

https://plotly.com/python/

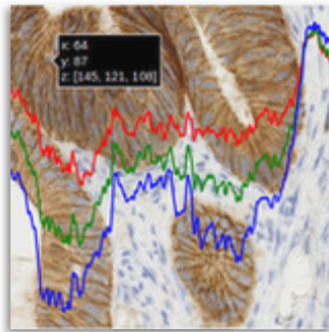# Python Plotly Graphing Library

## Scientific Charts
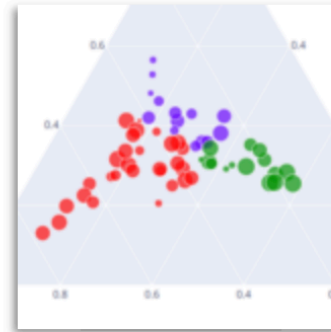
Contour Plots          Heatmaps          Imshow          Ternary Plots          Log Plots

## Financial Charts

Time Series and Date Axes          Candlestick Charts          Waterfall Charts          Funnel Chart          OHLC Charts

https://plotly.com/python/

# Python Plotly Graphing Library

## Maps

Mapbox Choropleth Maps

Lines on Mapbox

Filled Area on Maps

Bubble Maps

Mapbox Density Heatmap

## 3D Charts

3D Axes

3D Scatter Plots

3D Surface Plots

3D Subplots

3D Camera Controls

https://plotly.com/python/

# Python Plotly Graphing Library

## Subplots



Mixed Subplots



Map Subplots



Table and Chart Subplots



Figure Factory Subplots

## Jupyter Widgets Interaction



Plotly FigureWidget Overview



Jupyter Lab with FigureWidget



Interactive Data Analysis with FigureWidget ipywidgets



Click Events

https://plotly.com/python/

# Python Bokeh



https://bokeh.org/

# Python Altair

## Altair: Declarative Visualization in Python



Altair is a declarative statistical visualization library for Python, based on Vega and Vega-Lite, and the source is available on GitHub.

With Altair, you can spend more time understanding your data and its meaning. Altair's API is simple, friendly and consistent and built on top of the powerful Vega-Lite visualization grammar. This elegant simplicity produces beautiful and effective visualizations with a minimal amount of code.

## Getting Started

**GETTING STARTED**

Overview

Installation

Dependencies

Development Install

Basic Statistical Visualization

**GALLERY**

Example Gallery

**USER GUIDE**

Specifying Data in Altair

Encodings

Marks

Data Transformations

4.1.0

https://altair-viz.github.io/

69

# Iris flower data set

## setosa          versicolor      virginica

# Iris Classfication



Sepal

Petal

Versicolor

# iris.data

```
5.1,3.5,1.4,0.2,Iris-setosa
4.9,3.0,1.4,0.2,Iris-setosa
4.7,3.2,1.3,0.2,Iris-setosa
4.6,3.1,1.5,0.2,Iris-setosa
5.0,3.6,1.4,0.2,Iris-setosa
5.4,3.9,1.7,0.4,Iris-setosa
4.6,3.4,1.4,0.3,Iris-setosa
5.0,3.4,1.5,0.2,Iris-setosa
4.4,2.9,1.4,0.2,Iris-setosa
4.9,3.1,1.5,0.1,Iris-setosa
5.4,3.7,1.5,0.2,Iris-setosa
4.8,3.4,1.6,0.2,Iris-setosa
4.8,3.0,1.4,0.1,Iris-setosa
4.3,3.0,1.1,0.1,Iris-setosa
5.8,4.0,1.2,0.2,Iris-setosa
5.7,4.4,1.5,0.4,Iris-setosa
5.4,3.9,1.3,0.4,Iris-setosa
5.1,3.5,1.4,0.3,Iris-setosa
5.7,3.8,1.7,0.3,Iris-setosa
5.1,3.8,1.5,0.3,Iris-setosa
5.4,3.4,1.7,0.2,Iris-setosa
5.1,3.7,1.5,0.4,Iris-setosa
4.6,3.6,1.0,0.2,Iris-setosa
5.1,3.3,1.7,0.5,Iris-setosa
4.8,3.4,1.9,0.2,Iris-setosa
5.0,3.0,1.6,0.2,Iris-setosa
```

**setosa**



**virginica**



**versicolor**

# Iris Data Visualization

# Data Visualization in Google Colab

```
import seaborn as sns
sns.set(style="ticks", color_codes=True)
iris = sns.load_dataset("iris")
g = sns.pairplot(iris, hue="species")
```

```
import numpy as np
import pandas as pd
%matplotlib inline
import matplotlib.pyplot as plt
import seaborn as sns
from pandas.plotting import scatter_matrix
```

```python
# Import Libraries
import numpy as np
import pandas as pd
%matplotlib inline
import matplotlib.pyplot as plt
import seaborn as sns
from pandas.plotting import scatter_matrix
print('imported')
```

```
imported
```

```python
url = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"
names = ['sepal-length', 'sepal-width', 'petal-length', 'petal-width', 'class']
df = pd.read_csv(url, names=names)
print(df.head(10))
```

```python
# Load dataset
url = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"
names = ['sepal-length', 'sepal-width', 'petal-length', 'petal-width', 'class']
df = pd.read_csv(url, names=names)
print(df.head(10))
```

```
   sepal-length  sepal-width  petal-length  petal-width        class
0           5.1          3.5           1.4          0.2  Iris-setosa
1           4.9          3.0           1.4          0.2  Iris-setosa
2           4.7          3.2           1.3          0.2  Iris-setosa
3           4.6          3.1           1.5          0.2  Iris-setosa
4           5.0          3.6           1.4          0.2  Iris-setosa
5           5.4          3.9           1.7          0.4  Iris-setosa
6           4.6          3.4           1.4          0.3  Iris-setosa
7           5.0          3.4           1.5          0.2  Iris-setosa
8           4.4          2.9           1.4          0.2  Iris-setosa
9           4.9          3.1           1.5          0.1  Iris-setosa
```

# df.tail(10)

```python
print(df.tail(10))
```

|     | sepal-length | sepal-width | petal-length | petal-width | class |
|-----|-------------|-------------|--------------|-------------|-------|
| 140 | 6.7 | 3.1 | 5.6 | 2.4 | Iris-virginica |
| 141 | 6.9 | 3.1 | 5.1 | 2.3 | Iris-virginica |
| 142 | 5.8 | 2.7 | 5.1 | 1.9 | Iris-virginica |
| 143 | 6.8 | 3.2 | 5.9 | 2.3 | Iris-virginica |
| 144 | 6.7 | 3.3 | 5.7 | 2.5 | Iris-virginica |
| 145 | 6.7 | 3.0 | 5.2 | 2.3 | Iris-virginica |
| 146 | 6.3 | 2.5 | 5.0 | 1.9 | Iris-virginica |
| 147 | 6.5 | 3.0 | 5.2 | 2.0 | Iris-virginica |
| 148 | 6.2 | 3.4 | 5.4 | 2.3 | Iris-virginica |
| 149 | 5.9 | 3.0 | 5.1 | 1.8 | Iris-virginica |

# `df.describe()`

```
print(df.describe())
```

|       | sepal-length | sepal-width | petal-length | petal-width |
|-------|--------------|-------------|--------------|-------------|
| count | 150.000000   | 150.000000  | 150.000000   | 150.000000  |
| mean  | 5.843333     | 3.054000    | 3.758667     | 1.198667    |
| std   | 0.828066     | 0.433594    | 1.764420     | 0.763161    |
| min   | 4.300000     | 2.000000    | 1.000000     | 0.100000    |
| 25%   | 5.100000     | 2.800000    | 1.600000     | 0.300000    |
| 50%   | 5.800000     | 3.000000    | 4.350000     | 1.300000    |
| 75%   | 6.400000     | 3.300000    | 5.100000     | 1.800000    |
| max   | 7.900000     | 4.400000    | 6.900000     | 2.500000    |

# print(df.info())
# print(df.shape)

```
print(df.info())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 5 columns):
sepal-length     150 non-null float64
sepal-width      150 non-null float64
petal-length     150 non-null float64
petal-width      150 non-null float64
class            150 non-null object
dtypes: float64(4), object(1)
memory usage: 5.9+ KB
None
```

```
print(df.shape)
```
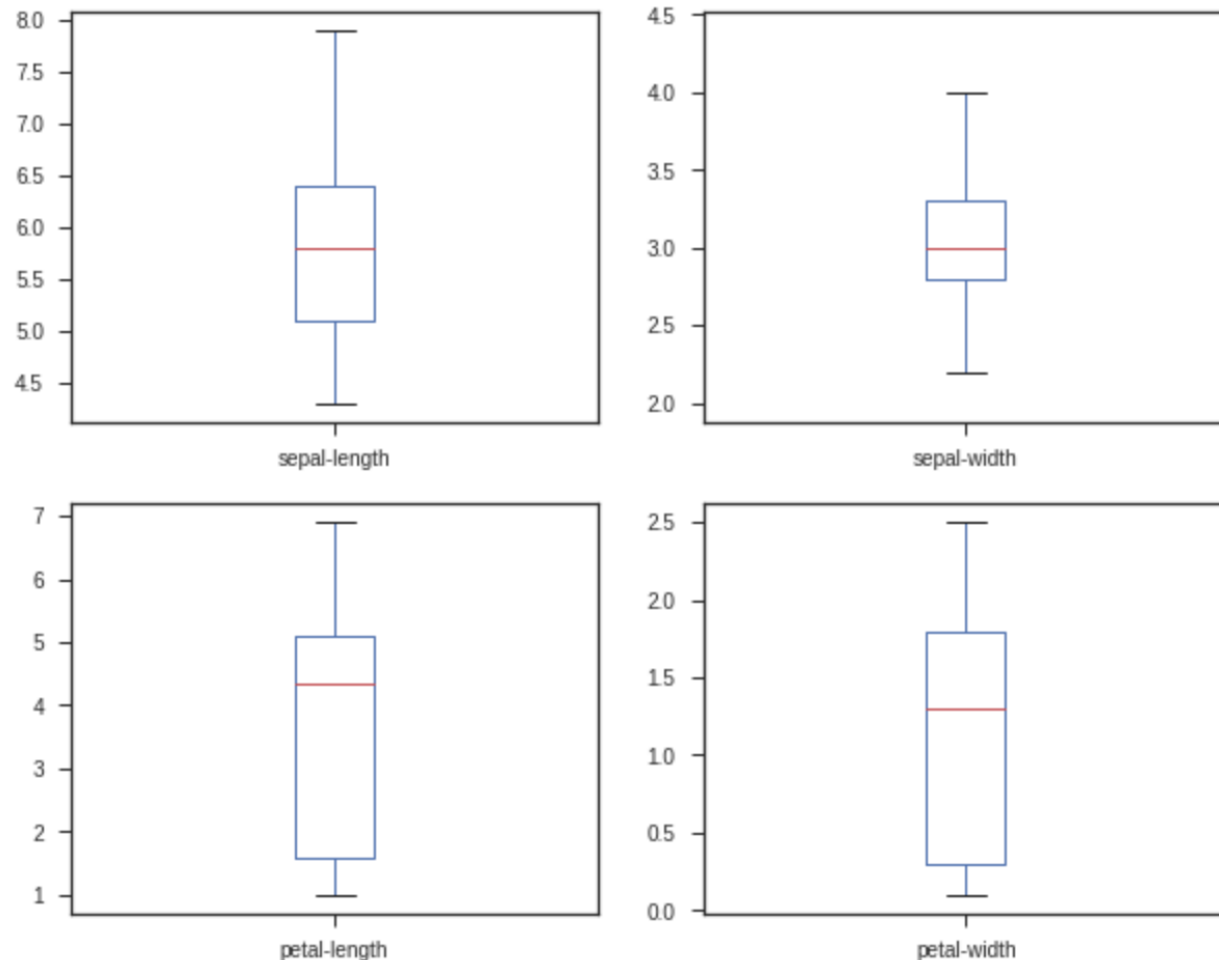
```
(150, 5)
```

# df.groupby('class').size()

```python
print(df.groupby('class').size())
```

```
class
Iris-setosa           50
Iris-versicolor       50
Iris-virginica        50
dtype: int64
```
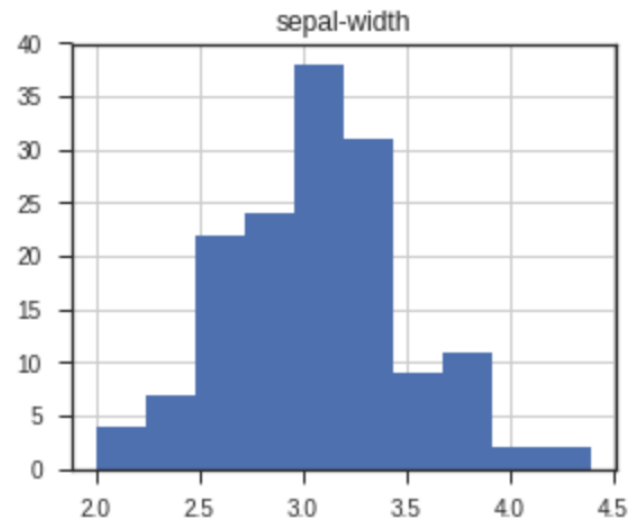
```
plt.rcParams["figure.figsize"] = (10,8)
df.plot(kind='box', subplots=True, layout=(2,2), sharex=False, sharey=False)
plt.show()
```
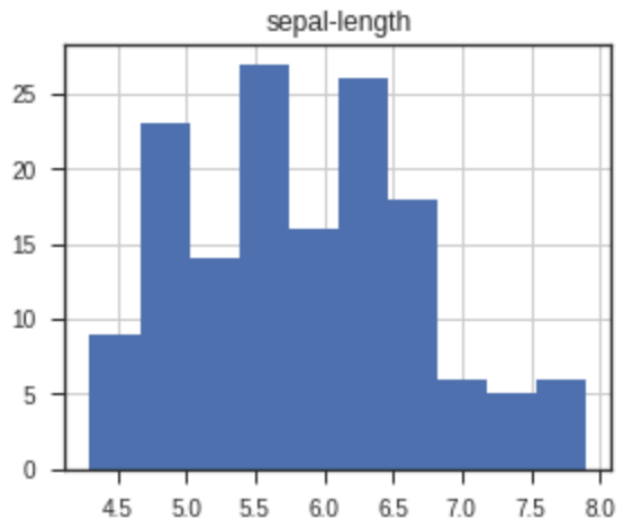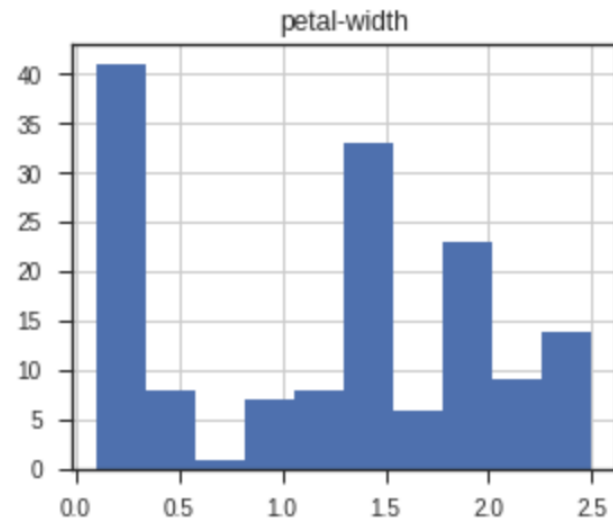
```
plt.rcParams["figure.figsize"] = (10,8)
df.plot(kind='box', subplots=True, layout=(2,2), sharex=False, sharey=False)
plt.show()
```
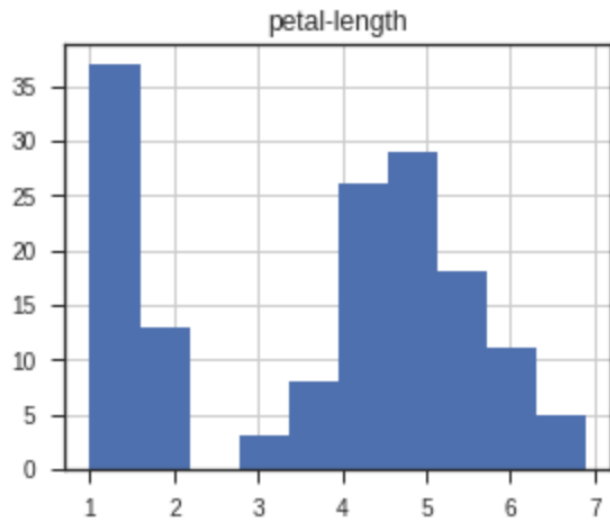
# df.hist()
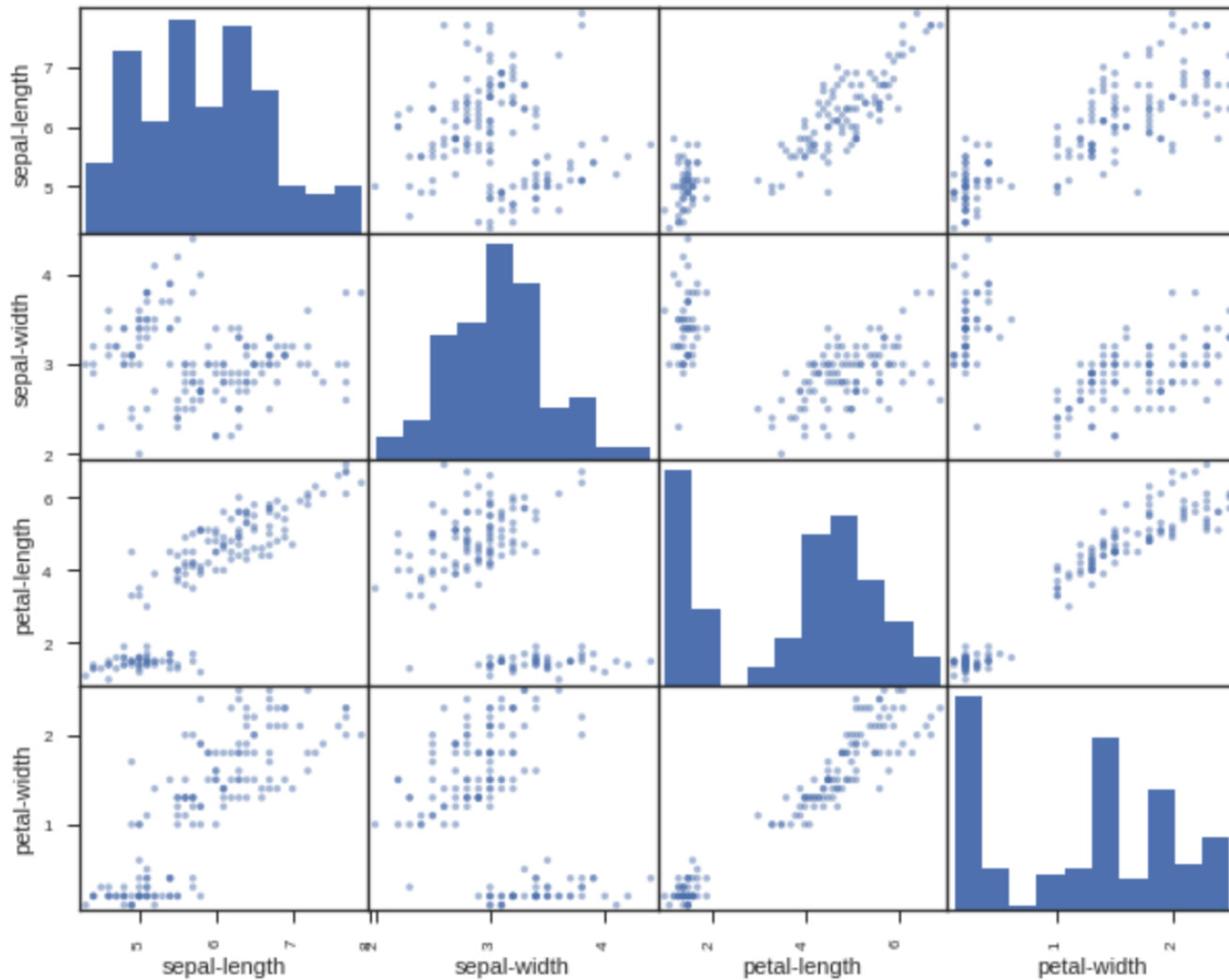# plt.show()

```
df.hist()
plt.show()
```

# scatter_matrix(df)
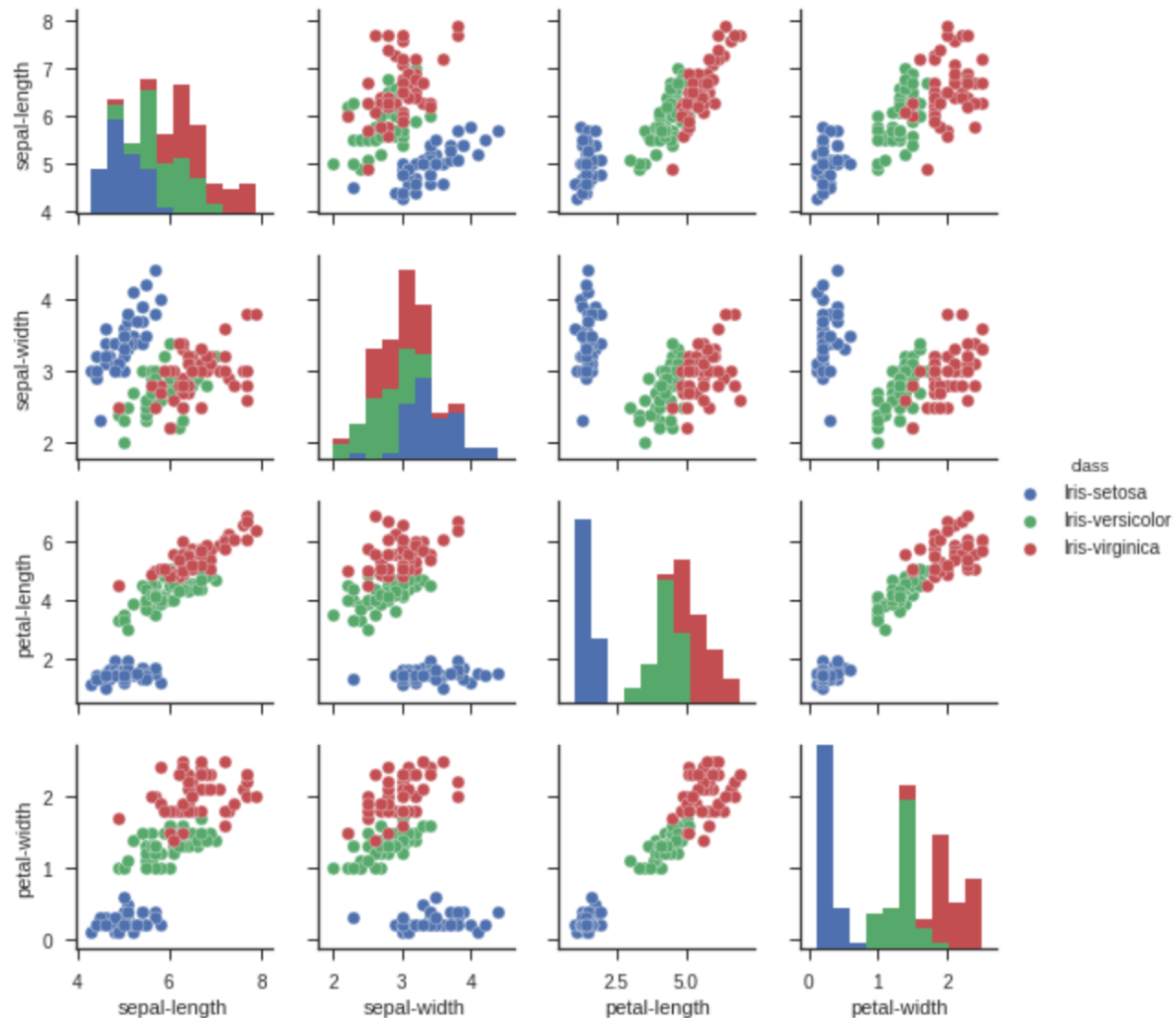# plt.show()

```
scatter_matrix(df)
plt.show()
```

# sns.pairplot(df, hue="class", size=2)

```
sns.pairplot(df, hue="class", size=2)
```

<seaborn.axisgrid.PairGrid at 0x7f1d21267390>

# Wes McKinney (2017), "Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython", 2nd Edition, O'Reilly Media.

Materials and IPython notebooks for "Python for Data Analysis" by Wes McKinney, published by O'Reilly Media

| | | | |
|---|---|---|---|
| ⊙ **52** commits | ⑂ **2** branches | ◇ **0** releases | ⚏ **6** contributors |

Branch: **2nd-edition ▾**    New pull request        Find file    **Clone or download ▾**

**betatim** committed with **wesm** Add requirements (#71)

| | |
|---|---|
| 📁 datasets | Add Kaggle titanic dataset |
| 📁 examples | Remove sex column from tips dataset |
| 📄 .gitignore | Add gitignore |
| 📄 COPYING | Use MIT license for code examples |
| 📄 README.md | Add launch in Azure Notebooks button (#70) |
| 📄 appa.ipynb | Make more cells markdown instead of raw |
| 📄 ch02.ipynb | Make more cells markdown instead of raw |
| 📄 ch03.ipynb | Make more cells markdown instead of raw |
| 📄 ch04.ipynb | Convert all notebooks to v4 format |
| 📄 ch05.ipynb | Make more cells markdown instead of raw |
| 📄 ch06.ipynb | Make more cells markdown instead of raw |
| 📄 ch07.ipynb | Convert all notebooks to v4 format |
| 📄 ch08.ipynb | Make more cells markdown instead of raw |
| 📄 ch09.ipynb | Make more cells markdown instead of raw |
| 📄 ch10.ipynb | Make more cells markdown instead of raw |

O'REILLY®
2nd Edition

# Python for Data Analysis

DATA WRANGLING WITH PANDAS, NUMPY, AND IPYTHON

powered by
jupyter

Wes McKinney

https://github.com/wesm/pydata-book

# Aurélien Géron (2019),
# Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems, 2nd Edition
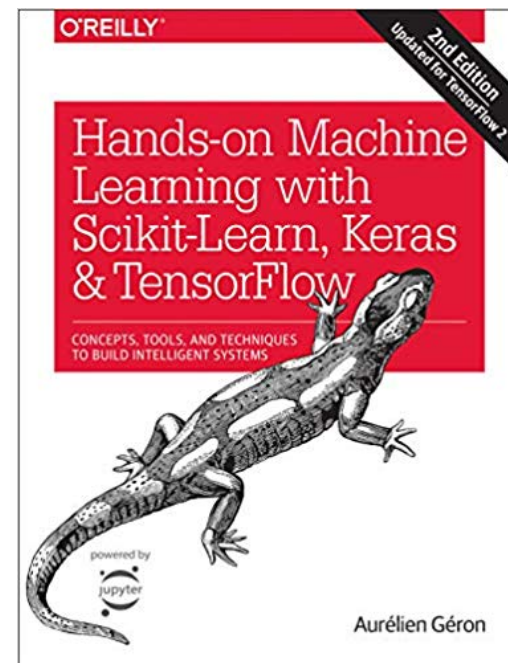# O'Reilly Media, 2019



https://github.com/ageron/handson-ml2

# Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow

**Notebooks**
1. The Machine Learning landscape
2. End-to-end Machine Learning project
3. Classification
4. Training Models
5. Support Vector Machines
6. Decision Trees
7. Ensemble Learning and Random Forests
8. Dimensionality Reduction
9. Unsupervised Learning Techniques
10. Artificial Neural Nets with Keras
11. Training Deep Neural Networks
12. Custom Models and Training with TensorFlow
13. Loading and Preprocessing Data
14. Deep Computer Vision Using Convolutional Neural Networks
15. Processing Sequences Using RNNs and CNNs
16. Natural Language Processing with RNNs and Attention
17. Representation Learning Using Autoencoders
18. Reinforcement Learning
19. Training and Deploying TensorFlow Models at Scale

https://github.com/ageron/handson-ml2

# Python in Google Colab (Python101)

https://colab.research.google.com/drive/1FEG6DnGvwfUbeo4zJ1zTunjMqf2RkCrT

# Summary

- **Data Science and Data Mining**

- **Discovering, Analyzing, Visualizing and Presenting Data with Python**
  - **Pandas**
  - **Matplotlib**
  - **Seaborn**
  - **Plotly**
  - **Bokeh, Altair**

# References

- EMC Education Services (2015),
  Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data, Wiley
- Ramesh Sharda, Dursun Delen, and Efraim Turban (2017), Business Intelligence, Analytics, and Data Science: A Managerial Perspective, 4th Edition, Pearson.
- Robert Layton (2017), Learning Data Mining with Python - Second Edition, Packt Publishing.
- Wes McKinney (2017), "Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython", 2nd Edition, O'Reilly Media.
- Aurélien Géron (2019), Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems, 2nd Edition, O'Reilly Media.
  https://github.com/wesm/pydata-book
- Pandas, http://pandas.pydata.org/
- Matplotlib, https://matplotlib.org/
- Seaborn, https://seaborn.pydata.org/
- Plotly, https://plotly.com/python/
- Bokeh, https://bokeh.org/
- Altair, https://altair-viz.github.io/
- Min-Yuh Day (2021), Python 101, https://tinyurl.com/aintpupython101