# **Open Source Data** Science and Al

### **SWE 406** 4<sup>th</sup> of July, 2025



**OLUWATODUNNI OBAFEMI** 22/10MSS013





# CONTENTS

- Open-source tools for data science and AI (Python ecosystem, **TensorFlow**, **PyTorch**)
- Collaborative data science with open tools
- Open datasets and their applications
- Ethical considerations in open AI development

# SECTION ONE: Open-Source Tools for Data Science and Al

#### What is open source?



• Fosters collaboration

Accelerates innovation

Democratizes access
 Enhances transparency and security

Why open source matters in AI?

# **Python: The Heart of Open Data Science**



Easy-to-learn syntax & vast ecosystem

Dominant language in Al and data science

#### Seamlessly integrates with AI frameworks

# **Core Libraries for Data Science**

#### NumPy

#### NUMERICAL COMPUTING

- An open source mathematical and scientific computing library for Python programming tasks.
- The name NumPy is shorthand for Numerical Pvthon.
- The NumPy library offers a collection of high-level mathematical functions including support for multidimensional arrays, masked arrays and matrices.

#### Pandas

#### DATA MANIPULATION

- A Python library used for working with data sets.
- It has functions for analyzing, cleaning, exploring, and manipulating data.
- The name "Pandas" has a reference to both "Panel Data", and "Python Data Analysis" and was created by Wes McKinney in 2008.

#### Matplotlib / Seaborn

DATA VISUALIZATION

- Matplotlib is a library in to generate visualizations like histograms, scatter and much more.
- Seaborn is a visualization visualizations that are statistically sophisticated.

Python that enables users plots, bar charts, pie charts

library that is built on top of Matplotlib. It provides data typically more aesthetic and

#### Scikit-learn

#### **MACHINE LEARNING ALGORITHMS**

- Scikit-learn's metrics enable thorough evaluation of machine learning models across different tasks and scenarios.
- Understanding these metrics helps in interpreting model performance, identifying potential areas for improvement and ultimately selecting or optimizing the bestperforming model for a specific problem.

# TensorFlow vs PyTorch



#### **TensorFlow**

Written in C++ and is, as a result, very fast and efficient.

Feature rich; TensorFlow can be used for training data as well as for inference.

Very good documentation; TensorFlow has many users and an big community which has led to strong documentation.

High popularity; TensorFlow has established itself as the most used ML library over a number of years now.

Many APIs available; TensorFlow is a library with a rich choice of easy to use APIs.

Supports JavaScript; TensorFlow supports JavaScript, C++ and Java in addition to Python.

For Mobile & IoT, inferences can be performed with TensorFlow Lite on mobile devices such as Android or iOS, as well as on Edge TPU or Raspberry Pi.



Written in Python making it more accessible and flattening the learning curve. However, the C++ core means PyTorch is still quite fast.

Very flexible; as data size can also be changed during data training.

Popular at research level; Pytorch was by far the most talked about ML library at CVPR, one of the most important computer vision conferences.

Rapid growth in popularity in both business and research use cases.

Many libraries available; PyTorch is composed of multiple libraries and platforms.

Python-based; PyTorch allows developers to write code in Python

PyTorch API; the PyTorch API is often preferred as it is better designed - plus TensorFlow has historically changed their API frequently.

# Additional Open Source Al Tools

AITOOL	USE
HuggingFace Transformers	Natural Language Processing state-of-the-art machine learnin and multimodal model, for both
Keras	High-level deep learning: Sim networks by providing a pytho
JAX	High Performance Machine Le numerical computing, particula features of NumPy with automa compilation.

**g:** acts as the model-definition framework for ng models in text, computer vision, audio, video, h inference and training.

plifies the processof building deep neural on interface for it

**earning:** a Python library for high-performance arly suited for machine learning, that combines atic differentiation and just-in-time (JIT)

### SECTION TWO:

### **Collaborative Data Science with Open Source Tools** The Need for Collaboration in Data Science

- Data projects are multidisciplinary
- Collaboration boosts quality and innovation
- Open tools remove barriers



### Jupyter Notebooks – A Data Scientist's Best Friend

- Interactive coding & output in one place
- Ideal for exploration, documentation, and sharing
- Supported by JupyterHub for teams

Jupyter Quickstart_Example (unsaved changes)       Image: Copy Bind         Sile       Sile	er link
File Eait View Insert Cell Kernel Widgets Help Trusted Pytho	on 3 O
Image: Constraint of the state of the	8 MB
In [25]: import matplotlib.pyplot as plt	
<pre>In [27]: # Temperature vs. Depth ax = plt.scatter(x=data['Temperature'],y=data['Depth'],c=data['Salinity' plt.gca().invert_yaxis(); # Flip the y-axis plt.colorbar(label='Salinity') plt.xlabel('Temperature (C)') plt.ylabel('Depth (m)');</pre>	]);
0 - 100 - - 35.6	
() 200 35.2 Lip 300 35.0	
400	



# Git & GitHub – Managing Code Together

• Version control to track changes

Enables peer review and contribution

• Fork, clone, pull request = smooth collaboration





### **Other Platforms for Open Collaboration**

• HuggingFace Hub Kaggle • Colab



BitBucket 

• GitLab



# SECTION THREE: OPEN DATASETS AND THEIR APPLICATIONS WHAT ARE OPEN DATASETS?

- Freely accessible data for public use
- Used for research, benchmarking, experimentation
- Often backed by governments or organizations



### **Top Sources for Open Datasets**

- Kaggle Datasets
   Google Dataset Search
  - UCI ML Repository





### World Bank WHO

#### Data.gov





## **Applications of Open Datasets**

#### **DRIVE RESEARCH AND POLICY TRAIN AND TEST AI MODELS**

#### SOLVE GLOBAL PROBLEMS **EDUCATIONAL USE AND SKILL** (CLIMATE, HEALTH, FINANCE) DEVELOPMENT



# **Example Case – ImageNet**



- 14 million+ labeled images
- Pioneered computer vision breakthroughs
- Basis for models like ResNet, VGG, YOLO



# SECTION FOUR: Ethical Considerations in Open Al Development

#### Why Ethics in Open AI?

- Al can amplify bias and injustice
- Open source = more eyes = more accountability
- Need for transparency and safeguards

#### **Responsibility vs Innovation**



# **Bias in Open Datasets**

- Historical data may reflect societal biases
- Models trained on biased data = biased predictions
- Need for diverse, representative datasets

Collect biased data





# **Transparency and Licensing**

- Open-source = visible decisions & algorithms
  - Choose licenses wisely (MIT, GPL, Apache)

License	Can be	Modifications	Can be re-	Contains special	Restrictiveness
	mixed with	can be taken	licensed	privileges for the	
	non-free	private and	by anyone	original copyright	
	software	not returned to		holder over user's	
		author		modifications	
General Public License (GPL)	No	No	No	No	Strong reciprocity
GNU Library General Purpose License	Yes	No	No	No	Standard reciprocity
Berkeley System Distribution (BSD)	Yes	Yes	No	No	Permissive
Netscape Public License	Yes	Yes	No	Yes	Standard reciprocity
Mozilla Public License <sup>2</sup>	Yes	Yes	No	No	Standard Reciprocity
Public Domain <sup>3</sup>	Yes	Yes	Yes	No	Permissive (Not OSS)

#### • Respect attribution, privacy, and terms

# **Preventing Misuse of Open Al**

- Dual-use dilemma: Good vs Harm
- Add documentation & use cases
- Consider ethical review before releas



# Future of Open-Source Al

- Global participation and innovation
- Ethical frameworks must grow alongside code
- Empowering the next generation of data scientists



# CONCLUSION

#### Utilization

#### **OPEN TOOLS + OPEN DATA**

#### Collaboration

Collaboration is the new superpower



#### Responsibility

### Build AI responsibly, build it together









# GUESTIONS?



