

Introduction to Machine Learning

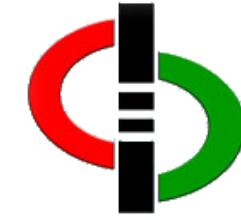
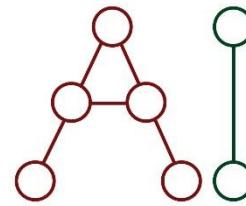
and Introduction to the AI 221 Course

Assoc. Prof. Karl Ezra Pilario, Ph.D.

Process Systems Engineering Laboratory

Department of Chemical Engineering

University of the Philippines Diliman



Outline

- What is Machine Learning?
 - Why only now?
 - Types of Learning Problems
- Intro to the Course (AI 221)
 - Course Delivery
 - Course Content
 - Course Requirements
 - Software

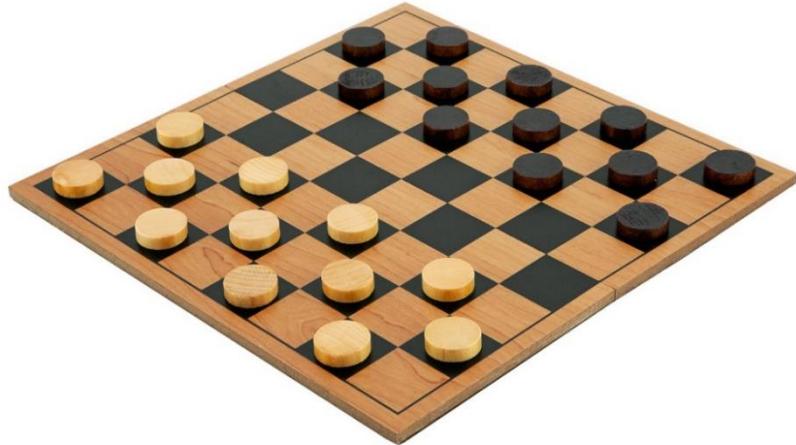
What is Machine Learning?



A field of study concerned with giving computers the *ability to learn* without being explicitly programmed.
(Arthur Samuel, 1959)

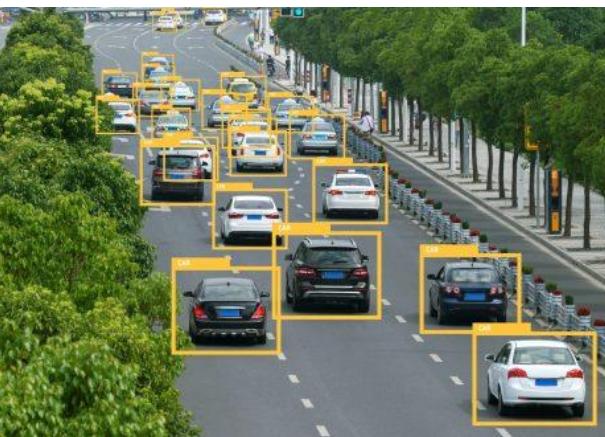


Arthur Samuel and the IBM 701 Computer



- Arthur Samuel was not a very good checkers player.
- But he was able to program a checkers bot that plays *better than a human*.
- The bot learned by playing millions of times against itself.
- In the process, it learns which moves lead to wins and which moves lead to losses.

0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9 9



K Explain quantum computing in simple terms

Quantum computing is a type of computing that uses quantum-mechanical phenomena, such as superposition and entanglement, to perform operations on data. Quantum computers are different from classical computers, which use bits to store and process information. Quantum computers use quantum bits, or qubits, which can represent a 0, a 1, or both at the same time.

One of the key differences between classical and quantum computers is that classical computers can only be in one state at a time, whereas quantum computers can exist in multiple states simultaneously. This property, known as superposition, allows quantum computers to perform certain types of calculations much faster than classical computers.

Another important difference is that quantum computers can exploit a phenomenon called

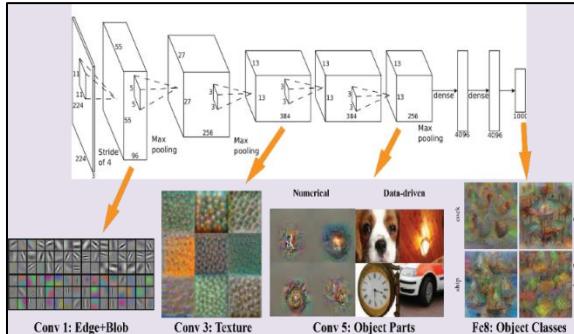


Typical ML Applications

- Filtering emails as spam / not spam
- Handwritten digits recognition
- Speech recognition, Natural Language Processing (NLP)
- Social media (Face recognition, News Feed Ranking, etc.)
- Image / Object Recognition, Image Segmentation
- Recommender systems (movies, products, videos, webpages, bookings)
- Finance (Stock market prediction, customer behaviour, etc.)
- Transportation (Self-driving cars, travel demand modelling)
- Healthcare (Early diagnostics, hospital demand forecasting)
- Bioinformatics (Protein folding and structure prediction, Gene function prediction, Biomedical image analysis)
- Chemometrics (GC-MS data analysis, drug discovery, compound classification, chemical property prediction)



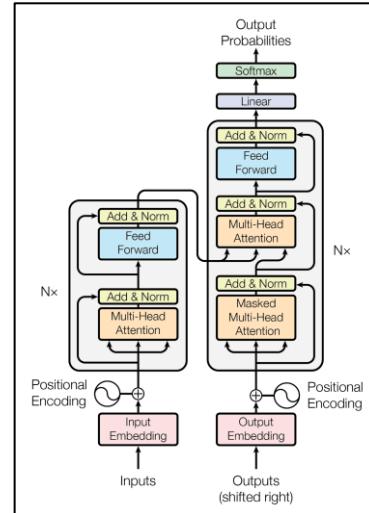
IBM Watson
Jeopardy, 2011



AlexNet
ImageNet Visual Recognition
Challenge, 2012



AlphaGo
Game of Go, 2016



Transformers
2017



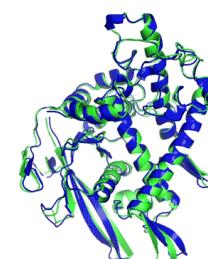
DALL-E
2021, 2022



IBM Deep Blue
Chess, 1997

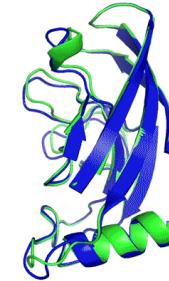


AlphaStar
StarCraft II, 2019



T1037 / 6vr4
90.7 GDT
(RNA polymerase domain)

- Experimental result
- Computational prediction



T1049 / 6y4f
93.3 GDT
(adhesin tip)

K Explain quantum computing in simple terms

Quantum computing is a type of computing that uses quantum-mechanical phenomena, such as superposition and entanglement, to perform operations on data. Quantum computers are different from classical computers, which use bits to store and process information. Quantum computers use quantum bits, or qubits, which can represent a 0, a 1, or both at the same time.

One of the key differences between classical and quantum computers is that classical computers can only be in one state at a time, whereas quantum computers can exist in multiple states simultaneously. This property, known as superposition, allows quantum computers to perform certain types of calculations much faster than classical computers.

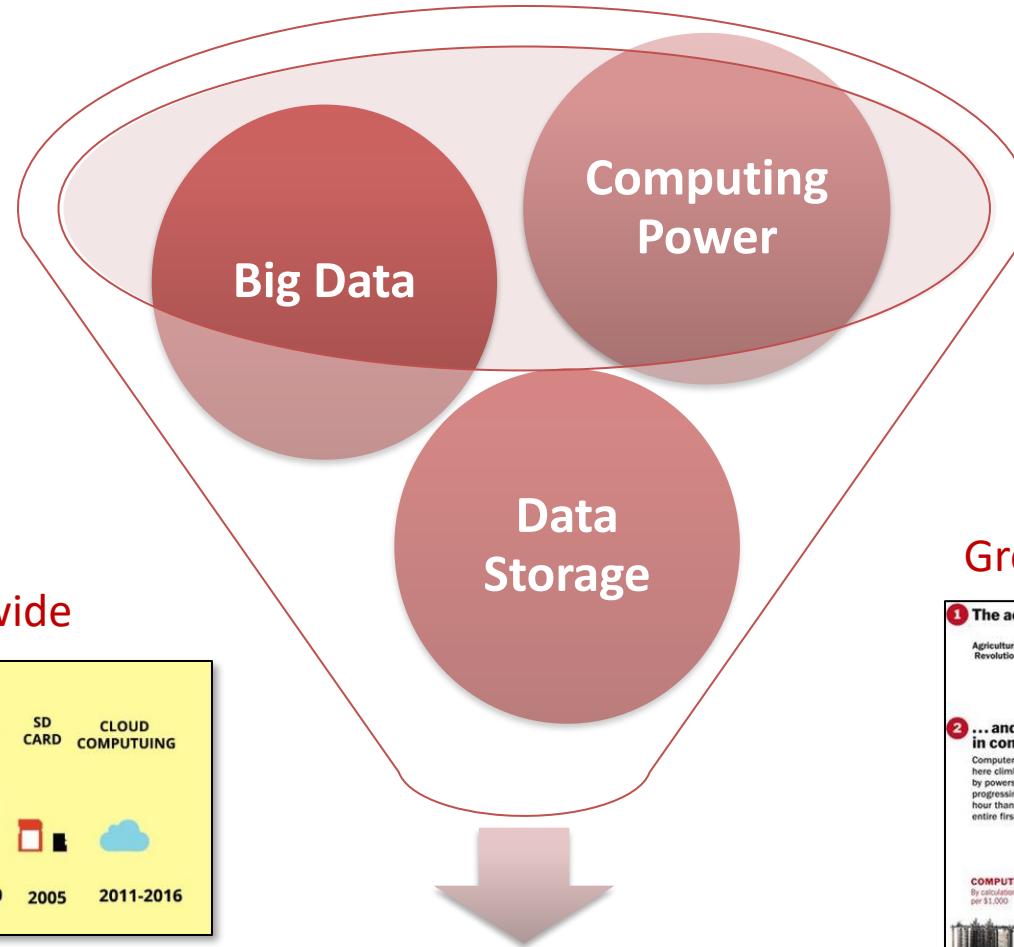
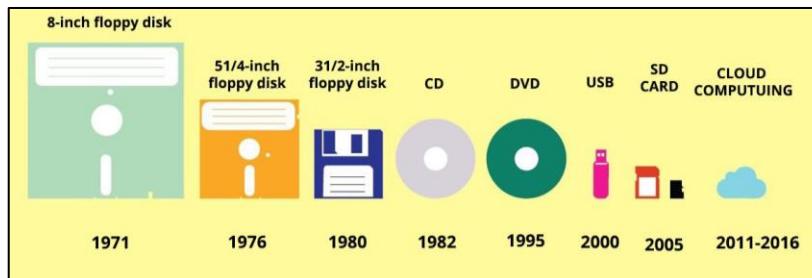
Another important difference is that quantum computers can exploit a phenomenon called entanglement, in which the state of one quantum particle can affect the state of another quantum particle, even if the two particles are separated by a large distance. This allows quantum computers to perform certain types of calculations in parallel, which

AlphaFold
Protein Structure Prediction,
2016, 2018

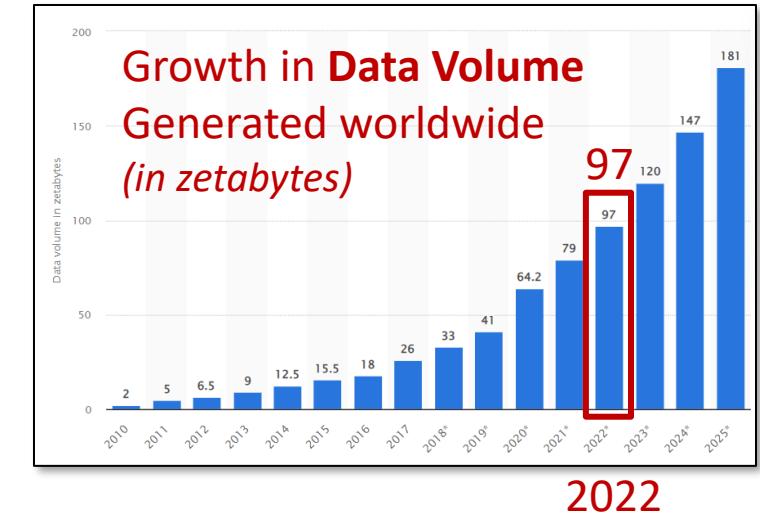
ChatGPT
2022

Machine Learning, Data Science, Data Analytics, ...why only now?

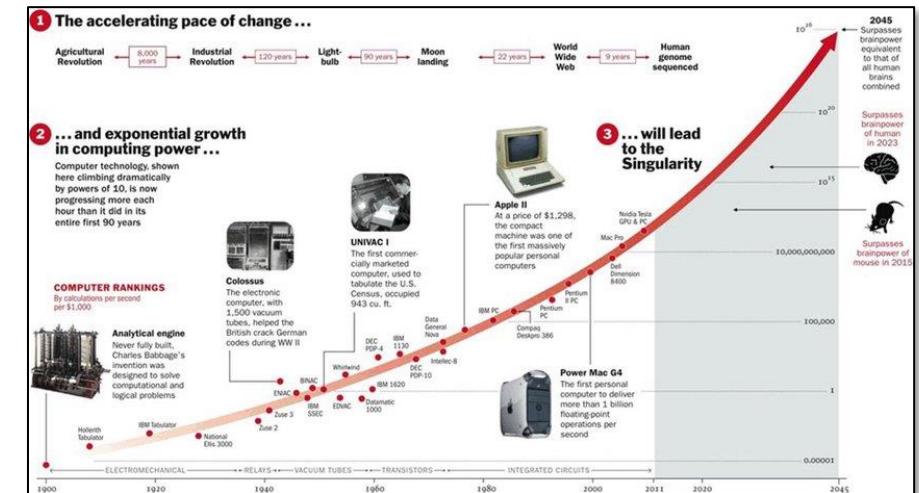
Growth in Data Storage worldwide



Machine Learning +
Practical Applications



Growth in Computing Power worldwide

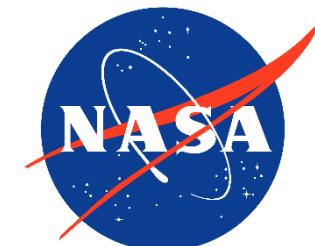


Machine Learning, Data Science, Data Analytics, ...why only now?

We are currently DROWNING¹ in data!

- There are about 1 trillion web pages.
- 1 hr of video is uploaded to Youtube every second.
- Human genomes have a length of 3.8×10^9 base pairs.
- Walmart handles more than 1 million transactions per hour.
- Etc...

Popular websites where we can get publicly available data:



Registry of Open Data on AWS

About

This registry exists to help people discover and share datasets that are available via AWS resources. See [recent additions](#) and [learn more about sharing data on AWS](#).

OpenML
A worldwide machine learning lab

Google Search public data

Public Data

Looking for other datasets? Find more with [Google Dataset](#)

Datasets Metrics Any data provider (45)

World Development Indicators World Bank This dataset contains the World Development Indicators (WDI).

THE WORLD BANK IBRD • IDA Data

¹ Venkatasubramanian (2009). DROWNING IN DATA: Informatics and Modeling Challenges in a Data-Rich Networked World. *AIChE Journal*.

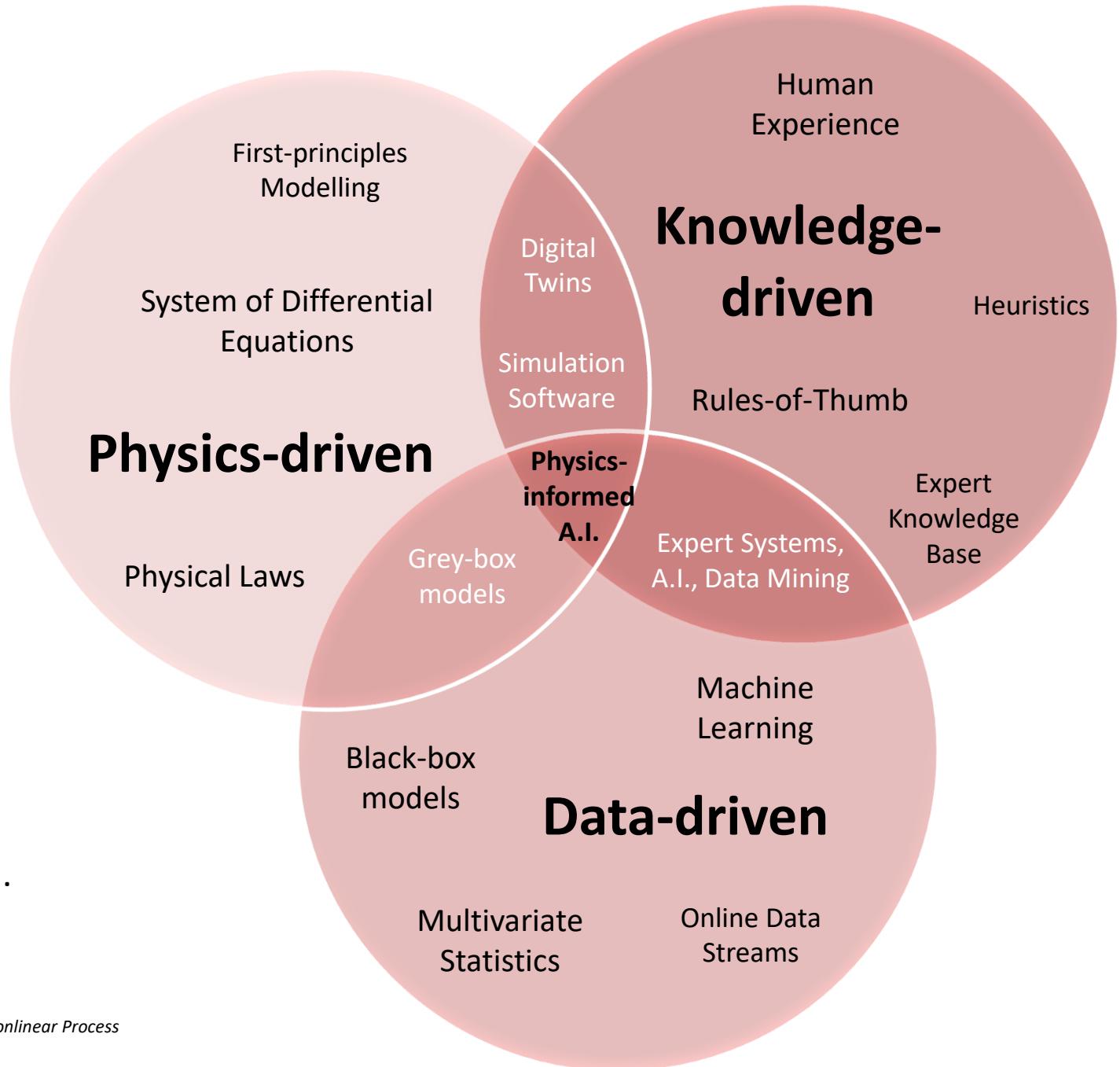
² Murphy (2012). Machine Learning: A Probabilistic Perspective. *MIT Press*.

Why use Machine Learning in your Industry?

Three approaches to engineering problems:

1. Physics-driven Methods
2. Knowledge-driven Methods
3. Data-driven Methods

Machine learning is a **data-driven approach**.



How to turn data into decisions?

Source: <https://iterationinsights.com/article/where-to-start-with-the-4-types-of-analytics/>

- Applying machine learning to your data is not enough.
- Don't just let your data speak, let it change the way you do things.
The goal is prescriptive analytics!
- Getting through each stage of analytics requires more and more effort, but also **more returns**.



Outline

- What is Machine Learning?
 - Why only now?
 - **Types of Learning Problems**
- Intro to the Course (AI 221)
 - Course Delivery
 - Course Content
 - Course Requirements
 - Software

Types of Learning Problems

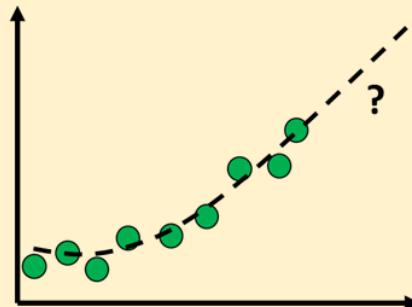
Supervised Learning

Learn a mapping or a function:

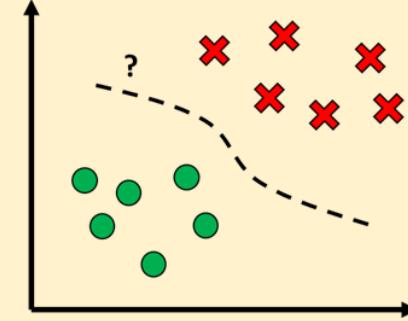
$$y = f(x)$$

from inputs (x) to outputs (y),
given a labelled set of input-output examples (● or ✗).

Regression



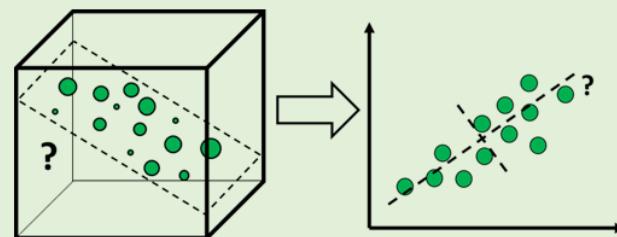
Classification



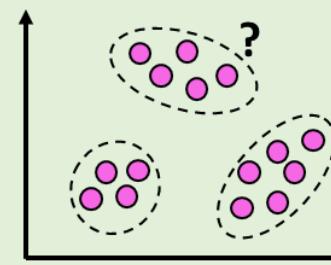
Unsupervised Learning

Discover *patterns or structure* from a data set (●) without any label information.

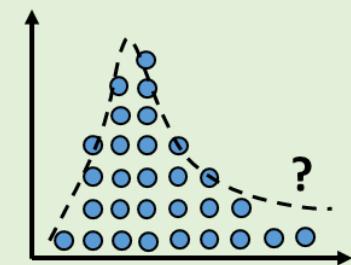
Dimensionality Reduction



Clustering



Density Estimation



Types of Learning Problems

A simple example...

Supervised Learning

These are images
of dogs.



Now, what is this
an image of?



These are
images of cars.



Unsupervised Learning

Here are some images...



Is there an image that does
not belong?

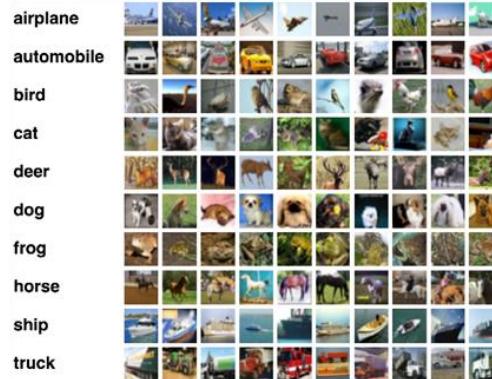
Are there images with similar
patterns?

Types of Learning Problems

Semi-Supervised Learning

Goal: Make a computer learn from both labelled and unlabelled data.

Labelled Data

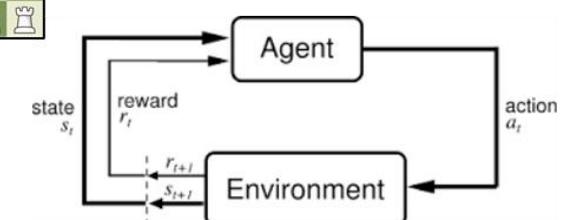


Unlabelled Data



Reinforcement Learning

Goal: Make a computer learn by letting it interact with the environment.



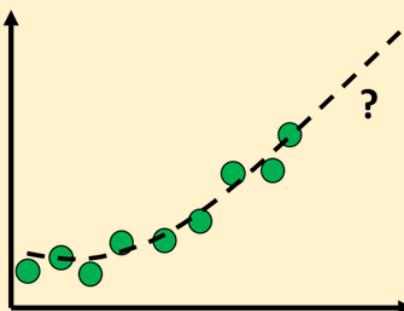
Supervised Learning

Learn a mapping or a function:

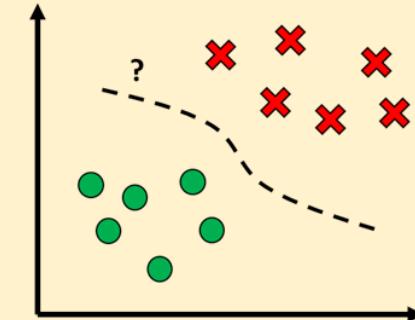
$$y = f(x)$$

from inputs (x) to outputs (y),
given a labelled set of input-output
examples (● or ✗).

Regression



Classification

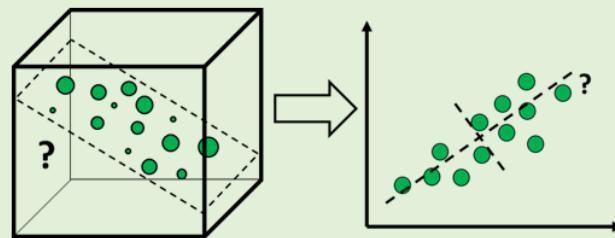


- **Given:** Training Data $\{x_i, y_i\}_{i=1,2,\dots,N}$
- Target y_i is a **continuous** variable.
- Examples:
 - Forecasting future stock price
 - Forecasting energy resources
 - Prediction of key performance indicators
 - Predicting the properties of molecules based on their structure
 - Predicting the environmental impact of pollutants
- **Given:** Training Data $\{x_i, y_i\}_{i=1,2,\dots,N}$
- Target y_i is a **categorical** variable.
- Examples:
 - Classifying objects in images
 - Classifying chest X-ray images into COVID positive/negative
 - Handwritten digits recognition
 - Filter e-mails into spam/not spam
 - Classify critical equipment as to healthy or faulty
 - Activity recognition from wearable devices

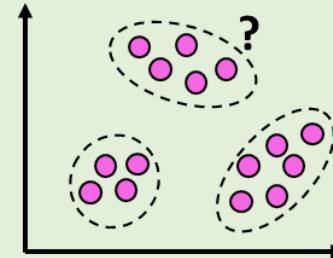
Unsupervised Learning

Discover *patterns or structure* from a data set (●) without any label information.

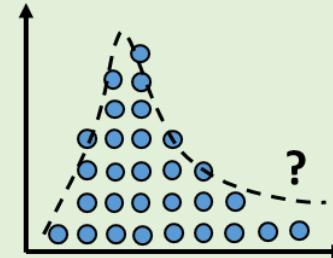
Dimensionality Reduction



Clustering



Density Estimation



Dimensionality Reduction

- **Given:** Data $\{x_i\}_{i=1,2,\dots,N}$
- **Reduce features** but retain the most important information from the original data.

Clustering

- **Given:** Data $\{x_i\}_{i=1,2,\dots,N}$
- **Group** similar data points together.

Density Estimation

- **Given:** Data $\{x_i\}_{i=1,2,\dots,N}$
- **Estimate** the distribution of the data.

Examples:

- Feature Engineering
- Image compression
- Filtering noise from signals
- Source separation in audio
- Data visualization

Examples:

- Customer segmentation
- Recommendation systems
- Identifying fake news
- Clustering documents, tweets, posts

Examples:

- Anomaly Detection
- Novelty Detection
- Generative Models
- Finding distribution modes
- Spatio-temporal analytics

Can you identify the type of learning problem?

Regression, Classification, Dimensionality Reduction, Clustering, Density Estimation

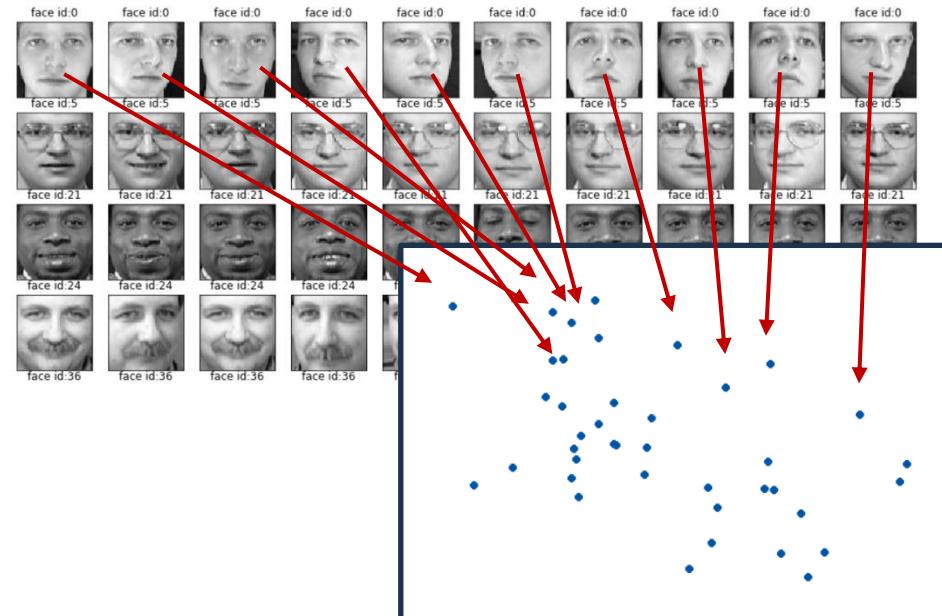
Example 1

Given the weight of the car, its model year, and horsepower, predict its mileage in miles per gallon (mpg).

car_weight	model_year	horsepower	mileage
1522 kg	2020	150	18 mpg
1930 kg	2017	185	16 mpg
1321 kg	2018	200	21 mpg
2128 kg	2019	168	?
2498 kg	2018	170	15 mpg
1882 kg	2021	155	17 mpg
1956 kg	2019	190	?
1672 kg	2017	182	18 mpg

Example 2

Given images of faces with varying poses and expressions, map each image onto a 2D point so that similar-looking images are closer together on the map.



Answer: Regression

Answer: Dimensionality Reduction

Can you identify the type of learning problem?

Regression, Classification, Dimensionality Reduction, Clustering, Density Estimation

Example 3

Given a tweet, predict whether the sentiment is positive, negative, or neutral.

Tweet	Sentiment
I'm in pain...	Negative
Manifesting a promotion this year!	Positive
It's 2AM. Who's awake?	Neutral
Heavy traffic at EDSA	Negative
Family dinner... So full!	Positive
Spoiler alert: RIP Tony Stark	?
Tesla sucks!	?
It's a boy!	Positive

Example 4

Given student grades in 5 subjects: Math, Chemistry, Physics, English, and Reading, group the students with similar competencies.

Student	Math	Chemistry	Physics	English	Reading
1	81	85	88	94	92
2	95	80	94	93	85
3	92	94	89	81	80
4	94	83	90	91	84
5	88	84	90	97	95
6	90	93	88	85	82
7	92	94	91	87	81
8	87	82	85	93	94

Answer: Classification

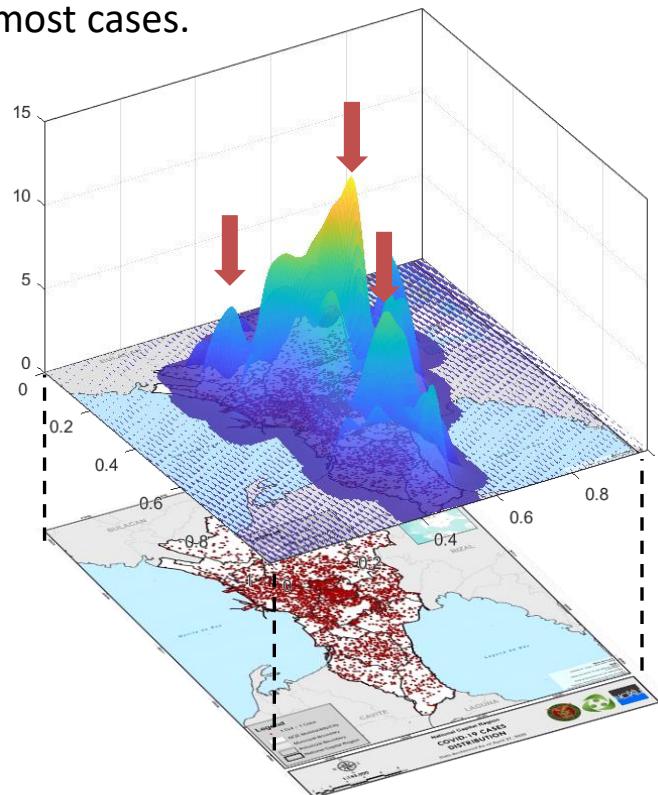
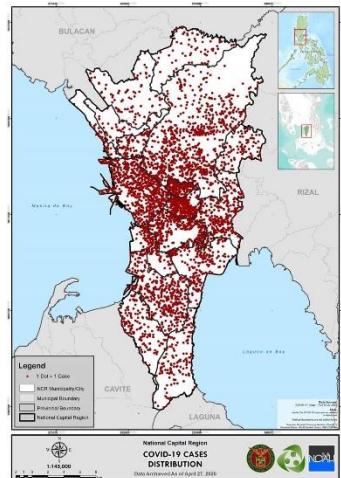
Answer: Clustering

Can you identify the type of learning problem?

Regression, Classification, Dimensionality Reduction, Clustering, Density Estimation

Example 5

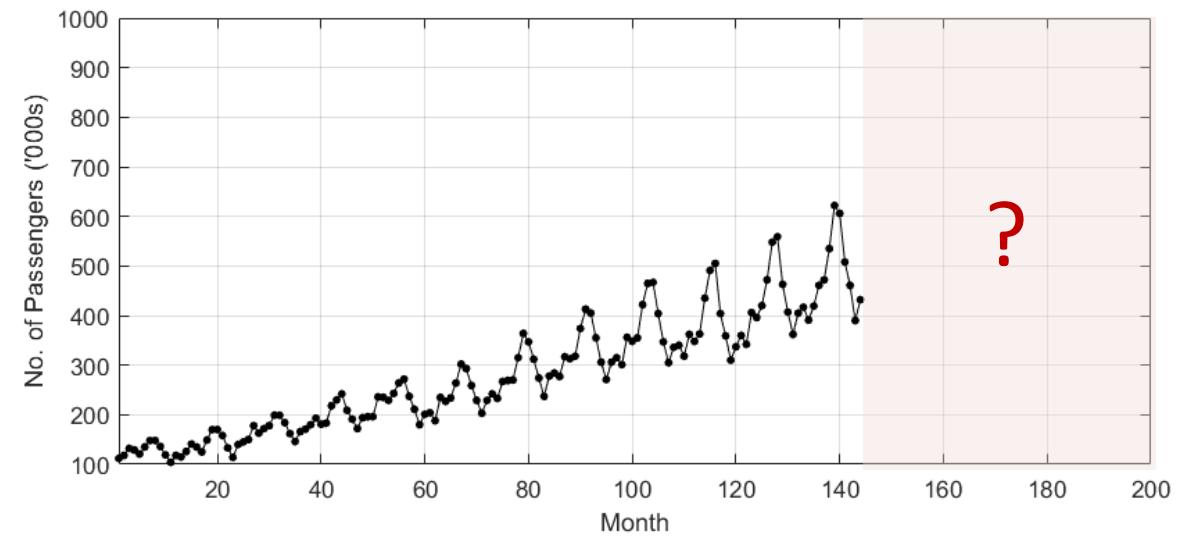
Given the spatial occurrence of Covid cases in Metro Manila, find the area with the most cases.



Answer: Density Estimation

Example 6

Given the number of airline passengers in the previous months, predict the number of passengers for the next few months.



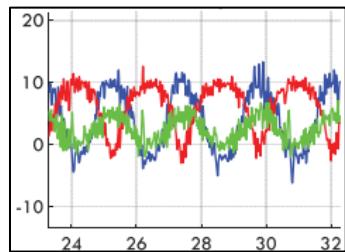
Answer: Regression

Can you identify the type of learning problem?

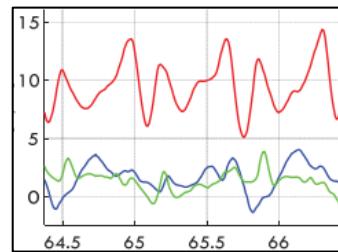
Regression, Classification, Dimensionality Reduction, Clustering, Density Estimation

Example 7

Given smartphone *accelerometer data* from a human doing exercise, predict the kind of exercise being done.

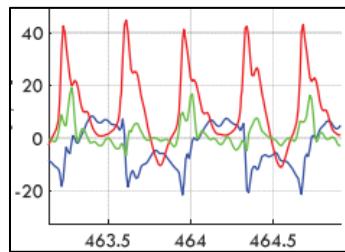


Push-up

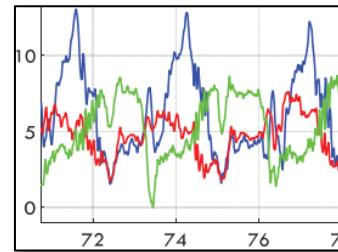


Example 8

Given the structural properties of alkane molecules, *map them onto 3D space based on their similarities, then predict their boiling points*.

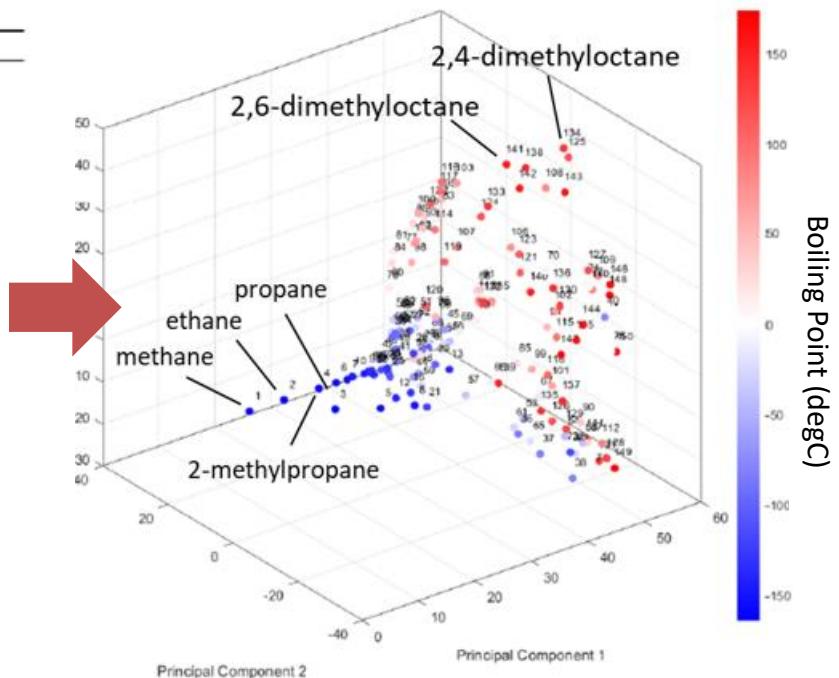


Running



Squats

No.	BP	Alkane
1	-164	methane
2	-88.6	ethane
3	-42.1	propane
4	-11.7	2-methylpropane
5	-0.5	butane
6	9.5	2,2-dimethylpropane
7	27.8	2-methylbutane
8	36.1	pentane
9	49.7	2,2-dimethylbutane
10	58	2,3-dimethylbutane
11	60.3	2-methylpentane
12	63.3	3-methylpentane
13	69	hexane
14	80.9	2,2,3-trimethylbutane
15	79.2	2,2-dimethylpentane
16	86.1	3,3-dimethylpentane
17	89.8	2,3-dimethylpentane
18	80.5	2,4-dimethylpentane
19	90	2-methylhexane
20	92	3-methylhexane
21	92.4	2-ethylhexane

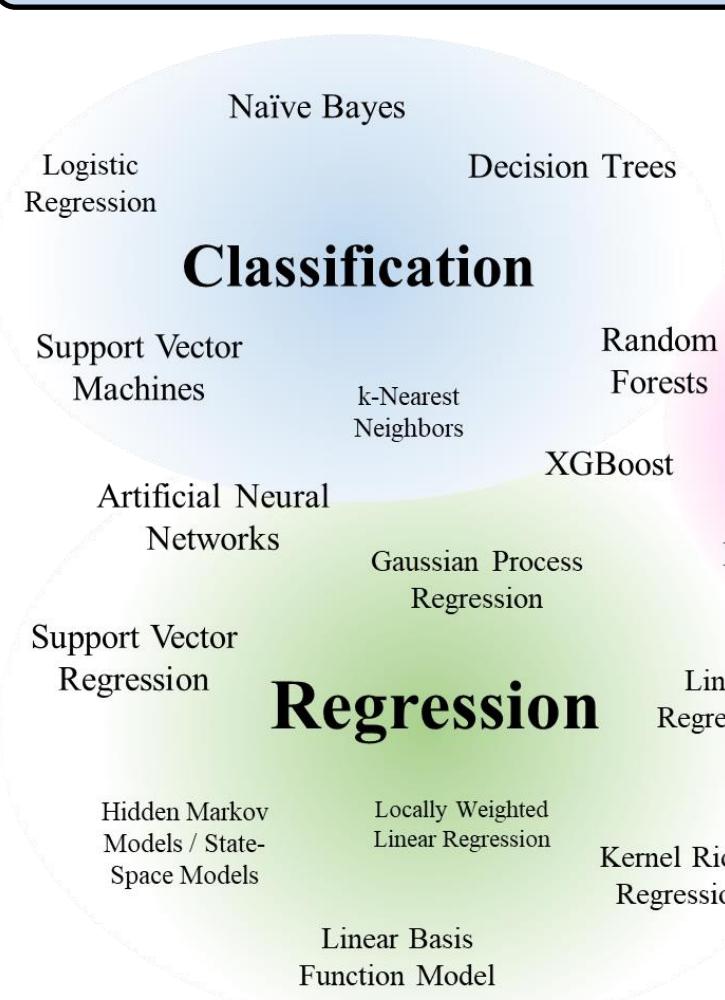


Answer: Classification

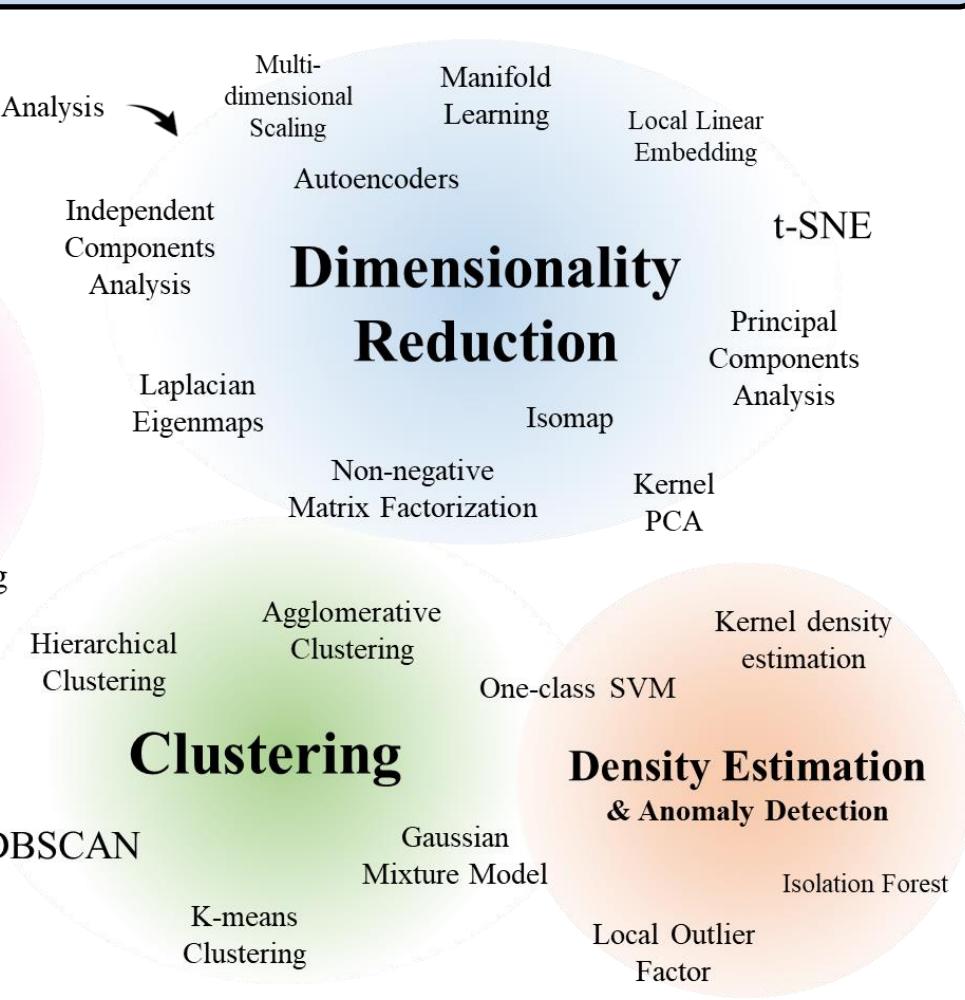
Answer: Dimensionality Reduction + Regression

Machine Learning Methods

Supervised Learning



Unsupervised Learning



Outline

- What is Machine Learning?
 - Why only now?
 - Types of Learning Problems
- **Intro to the Course (AI 221)**
 - Course Delivery
 - Course Content
 - Course Requirements
 - Software

Introduction to the Course

COURSE NUMBER:

AI 221

COURSE TITLE:

Classical Machine Learning

COURSE DESCRIPTION:

Linear Models. Kernel Methods. Neural Networks. Trees. Clustering. Dimensionality Reduction. Feature Engineering. Density Estimation. Ensemble Learning. Gaussian Processes. Bayesian Methods. Hyperparameter Search. AutoML. Explainability.

COURSE CREDIT:

3 units

3.0 hours/week

COURSE LMS*:

UVLE Course Page: AI 221 [WZZQ]

Github:

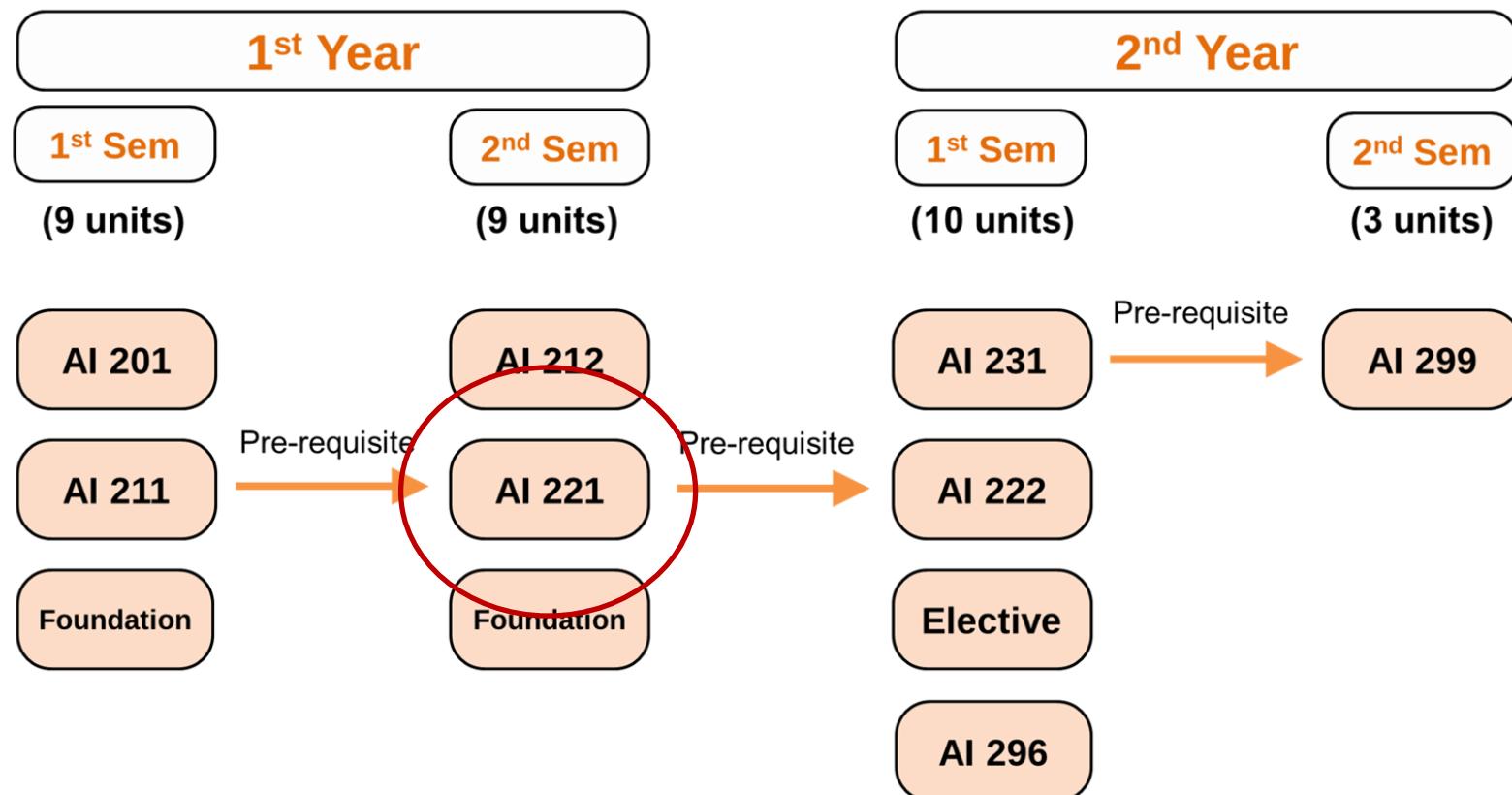
<https://github.com/kspilario/AI221>

*LMS = Learning Management System

Introduction to the Course

MEngg in Artificial Intelligence

Total Units: 31 Units



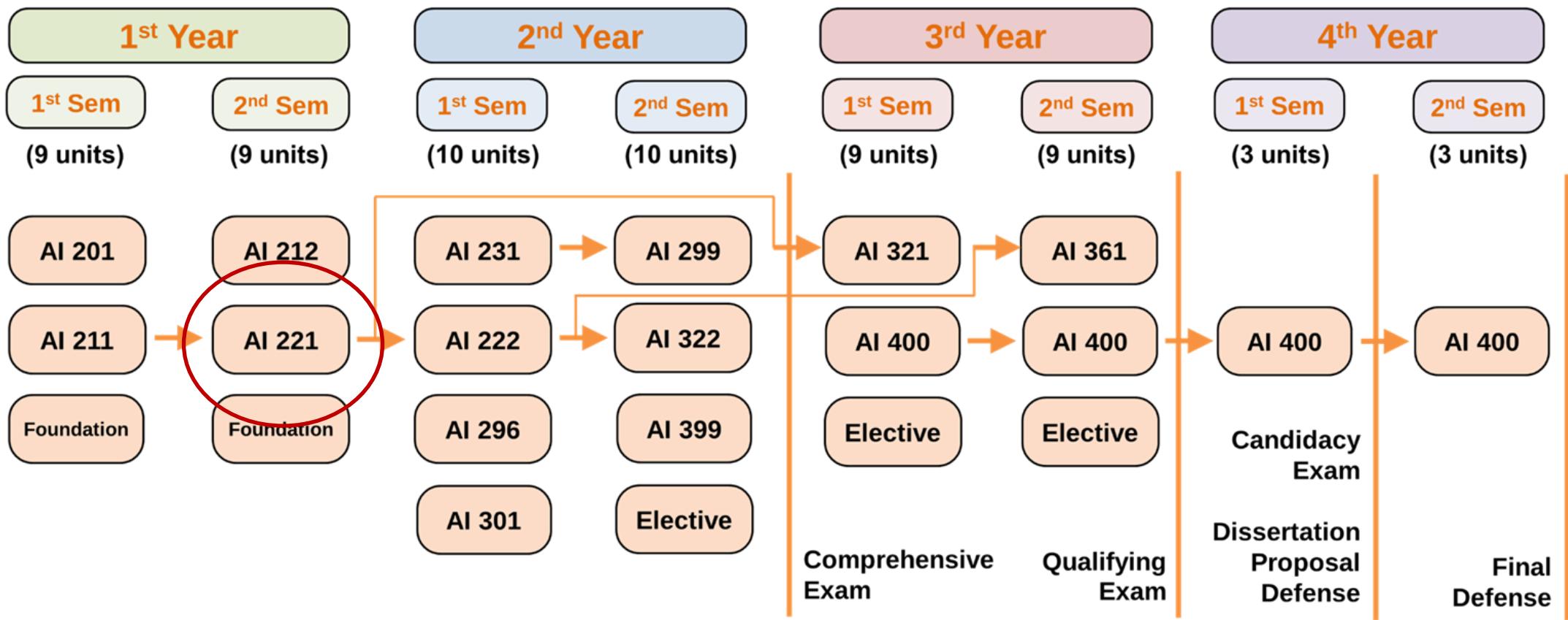
Introduction to the Course

PhD in Artificial Intelligence

Option A

Total Units: 62 Units

Requirement: 2 Publications



AI 221 Course Delivery

- Meeting:** Every Tuesday, face-to-face, Room A301, Chemical Engineering Building.

Sunday	Monday	Tuesday	Wednesday	Thursday	Friday	Saturday
			6-9 PM			

Course Requirements:	Requirement	% of Final Grade	Mode
	<ul style="list-style-type: none">Team Project<ul style="list-style-type: none">Oral Presentation (40%)Written Report (50%)Graphical Abstract (10%)	40%	By groups of 1 to 3 members only, Face-to-face
	<ul style="list-style-type: none">Machine Exercises	40%	By group (1-3), Take-home
	<ul style="list-style-type: none">Journal Critique	20%	Individual, Take-home

- Grading System:**

[92,100]	[88,92)	[84,88)	[80,84)	[76,80)	[72,76)	[68,72)	[64,68)	[60,64)	[0,60)
1.00	1.25	1.50	1.75	2.00	2.25	2.50	2.75	3.00	5.00

AI 221 Course Content

Feb 7 **Week 1.** Introduction to Machine Learning

Feb 14 **Week 2.** Exploratory Data Analysis

Feb 21 **Week 3.** Linear and Logistic Regression

Feb 28 **Week 4.** Support Vector Machines and Kernel Methods

Mar 6 **Week 5.** Cross-validation and Hyper-parameter Optimization

Mar 13 **Week 6.** Linear Dimensionality Reduction + Discriminant Analysis

Mar 20 **Week 7.** Nonlinear Dimensionality Reduction

----- Reading Break + Lenten Break -----

Apr 10 **Week 8.** Clustering, Density Estimation, and Anomaly Detection

Apr 17 **Week 9.** Trees, Weak Learners, and Ensemble Learning (Boosting, Bagging, Stacking)

Apr 24 **Week 10.** Neural Networks for Classification, Regression, and Dim. Reduction

May 8 **Week 11.** Gaussian Processes and Bayesian Optimization

May 15 **Week 12.** AutoML and ML Explainability

May 22 **Week 13. Team Project Presentation**

AI 221 Course Requirements

Team Project (40%)

- A team should have **at most 3 members** only.
- Aims:
 - Find a **problem + data set** that requires an ML solution.
 - Solve the problem using the **ML methods** discussed in class.
 - **Present** your results to the class.
- **NO** two teams should have the same problem.
- Grading and deadlines:
 - Oral Presentation (40%) – **May 22, 2024**
 - Written Report (50%) – **May 29, 2024**
 - Graphical Abstract (10%) – **May 29, 2024**

Machine Exercises (40%)

- Mode: **Groups of 1-3 members, take-home**
- To be given every lecture.
- Submission deadline is **2 weeks** after release.

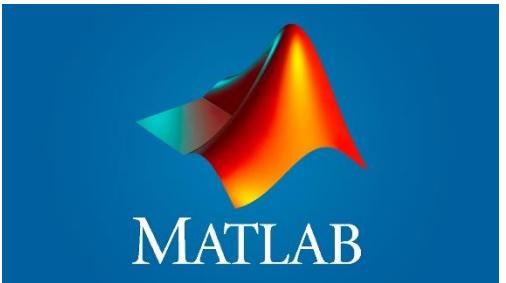
Journal Critique (20%)

- Mode: **Individual, take-home**
- Find a paper from a reputable journal or conference proceedings related to your field.
 - Should at least have an impact factor.
 - Should be published in the **last 5 years**.
- **Send me** the paper for approval first, then I will send guide questions for you to answer.
- Deadline for Paper approval: **Jun 7, 2024**
- Deadline for Critique submission: **Jun 14, 2024**

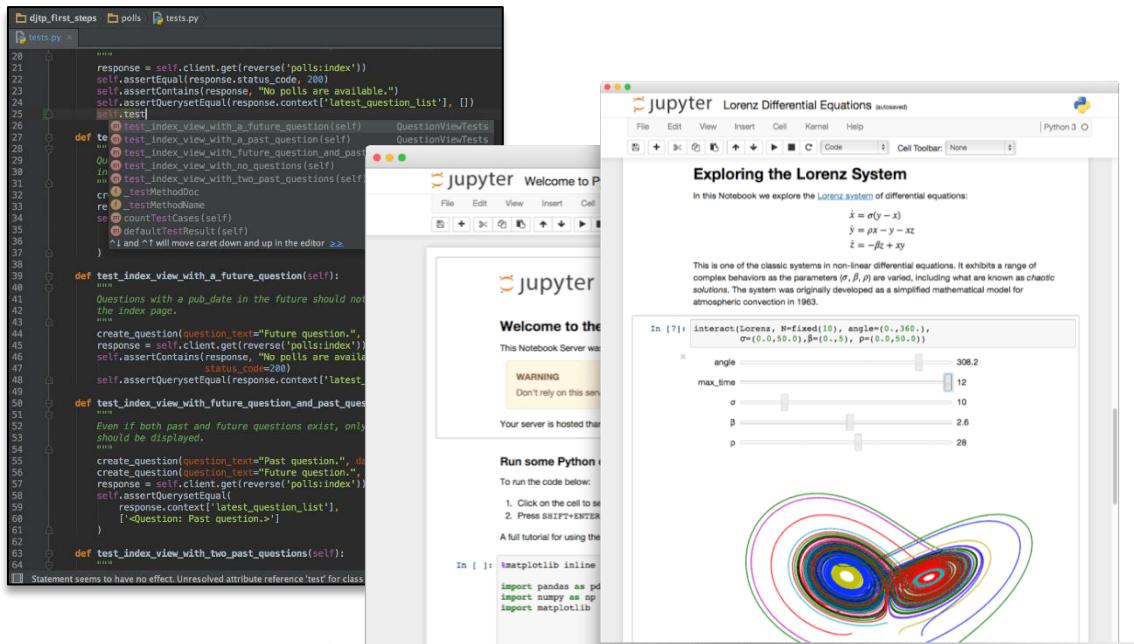
AI 221 Required Software



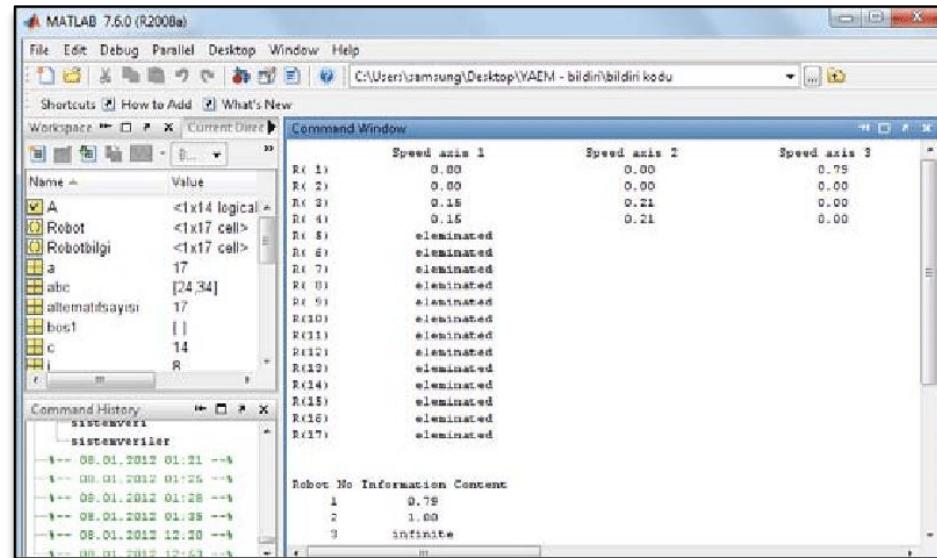
or



- Anaconda >> Spyder
- Anaconda >> JupyterLab
- Google Colab
- Jupyter Notebook
- PyCharm
- JupyterLite
- Microsoft Excel (**New**)



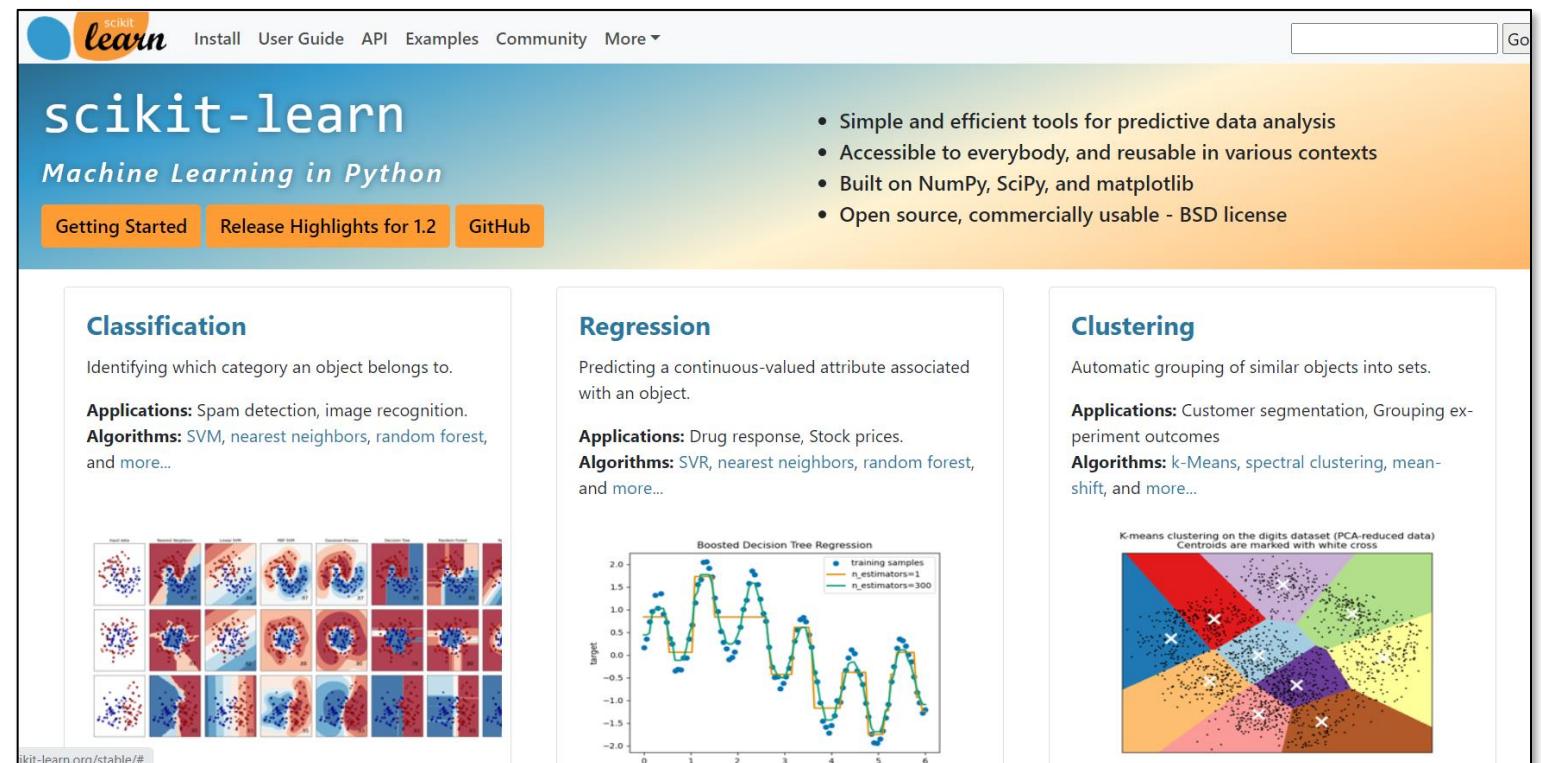
- You can download MATLAB by logging in to www.mathworks.com
 - Use your UP credentials!
- You can also use MATLAB online.



AI 221 Required Software

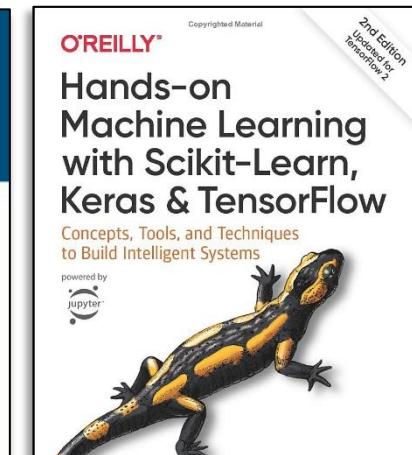
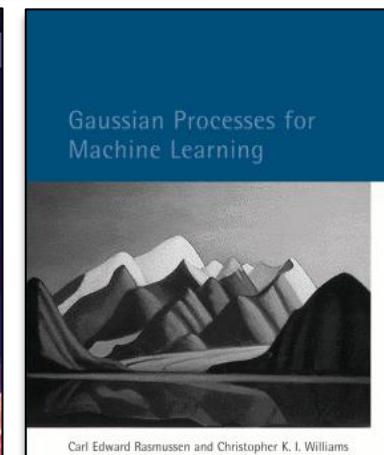
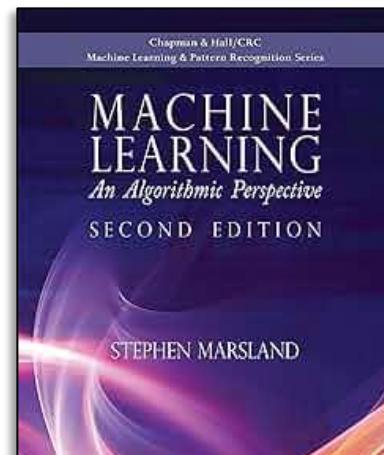
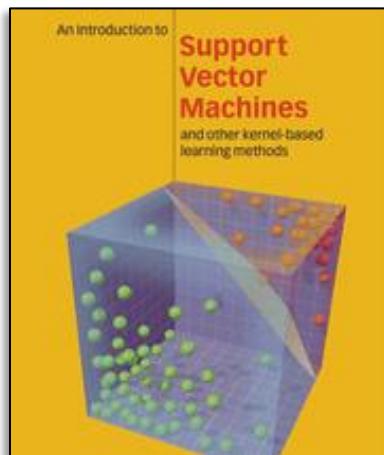
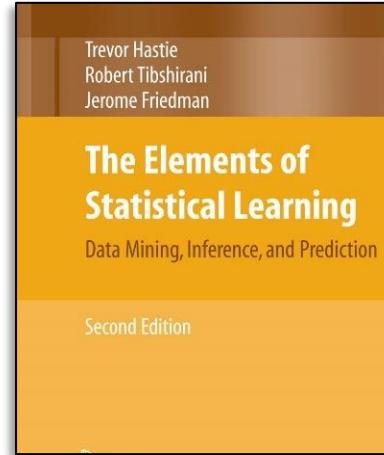
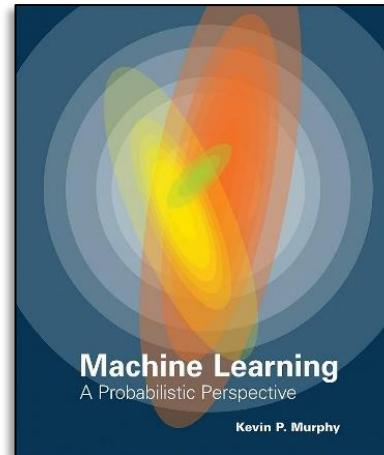
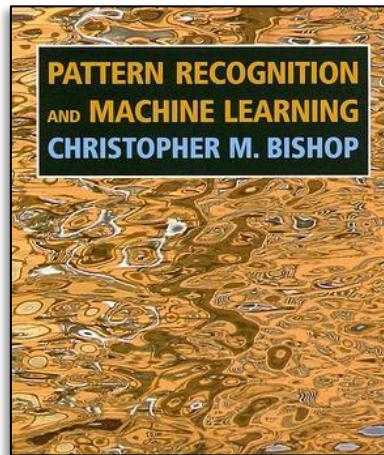


- Python 3
 - <https://www.python.org/>
- Numpy
 - <http://www.numpy.org/>
- Scikit-Learn
 - <https://scikit-learn.org/>
- Jupyter Lab
 - <https://jupyter.org/try-jupyter/lab/>
 - <https://nbviewer.org/>
- MS Excel

A screenshot of the scikit-learn homepage. The header features the scikit-learn logo and navigation links for Install, User Guide, API, Examples, Community, and More. Below the header, the title "scikit-learn" and subtitle "Machine Learning in Python" are displayed. A navigation bar includes links for Getting Started, Release Highlights for 1.2, and GitHub. To the right, a list of features is presented: "Simple and efficient tools for predictive data analysis", "Accessible to everybody, and reusable in various contexts", "Built on NumPy, SciPy, and matplotlib", and "Open source, commercially usable - BSD license". The main content area is divided into three sections: Classification, Regression, and Clustering. Each section contains a brief description, applications, algorithms, and associated visualizations. The Classification section shows a grid of 9x3 plots for various classification models. The Regression section shows a line plot titled "Boosted Decision Tree Regression" comparing training samples, n_estimators=1, and n_estimators=300. The Clustering section shows a scatter plot titled "K-means clustering on the digits dataset (PCA-reduced data)" with centroids marked by white crosses.

AI 221 References

- Bishop (2006). *Pattern Recognition and Machine Learning*. Springer.
- Murphy, Kevin (2012). *Machine Learning: A Probabilistic Perspective*. MIT Press.
- Hastie et al. (2008). *The Elements of Statistical Learning*. 2nd Ed. Springer.
- Cristianini & Shawe-Taylor (2000). *An Introduction to Support Vector Machines and other kernel-based learning methods*. Cambridge University Press.
- Marsland, Stephen (2014). *Machine Learning: An Algorithmic Perspective*. Chapman and Hall. 2nd Ed.
- Rasmussen and Williams (2006). *Gaussian Processes for Machine Learning*. MIT Press.
<https://gaussianprocess.org/gpml/>
- Geron, Aurelien (2019). *Hands-on Machine Learning with Scikit-Learn, Keras & TensorFlow*. O'Reilly Media.
- Journals and Conference Proceedings
- Python API, Sci-kit learn API: <https://scikit-learn.org/stable/modules/classes.html>
- Online Courses, Youtube Videos, etc.



AI 221 Course Instructor



Current Position

Karl Ezra S. Pilario

Associate Professor

Department of Chemical Engineering
University of the Philippines, Diliman

- Process Dynamics & Control
- Programming in MATLAB, Python, Aspen HYSYS
- Numerical Methods in Engineering
- Plant Design and Research
- Machine Learning and Artificial Intelligence

Education

- **Bachelor's Degree:**

Chemical Engineering, SCL (**2012**)
University of the Philippines Diliman

- **Master's Degree:**

Chemical Engineering (**2015**)
University of the Philippines Diliman

- **PhD Degree:**

PhD Energy and Power (**2020**)
Cranfield University, United Kingdom



University of
the Philippines,
Diliman

Cranfield University, U.K.



Research Lab

Head, Process Systems Engineering Laboratory (PSEL)

Department of Chemical Engineering
University of the Philippines - Diliman

Research Interests

- Process Data Analytics
- Process Systems Engineering
- Industrial Process Monitoring and Predictive Maintenance
- Machine Learning for Energy, Water, and Environment
- Cheminformatics and Materials Informatics

Outline

- What is Machine Learning?
 - Why only now?
 - Types of Learning Problems
- Intro to the Course (AI 221)
 - Course Delivery
 - Course Content
 - Course Requirements
 - Software