

ML with Azure Scholarship

Lesson 1: Introduction to ML

1) What is ML?

→ It is a data science technique used to extract pattern from data, allowing computers to identify data & forecast future outcomes, behaviour & trend.

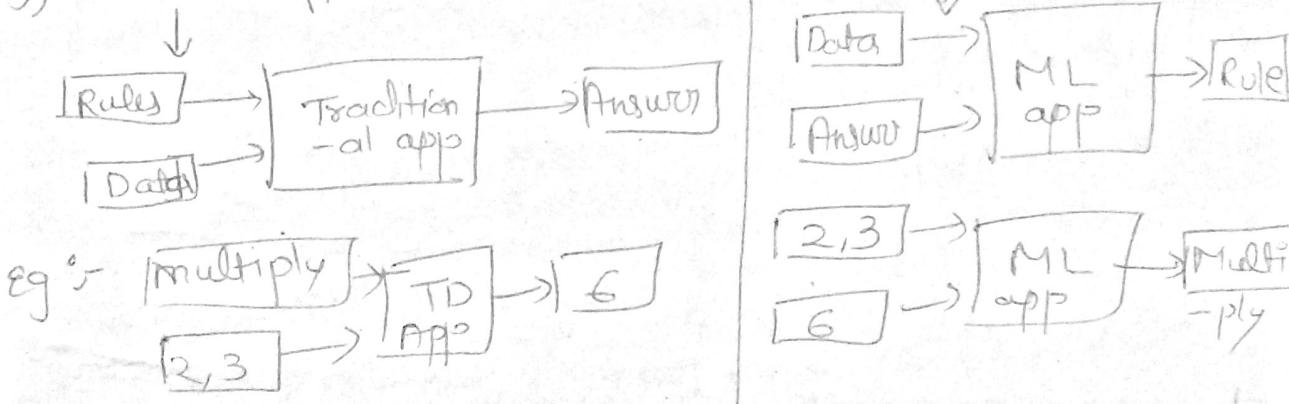
2) Application :- It is used in

- 1) NLP
- 2) Computer Vision
- 3) Data Analysis

4) Decision Making

- i) Automate the recognition of disease
- ii) Recommended next best action for individual customer
- iii) Chatbot
- iv) Identify the next best action for the customer

3) Tradition approach VS ML approach



4) Diff Between

AI

→ A broad term that refers to computers thinking more like human

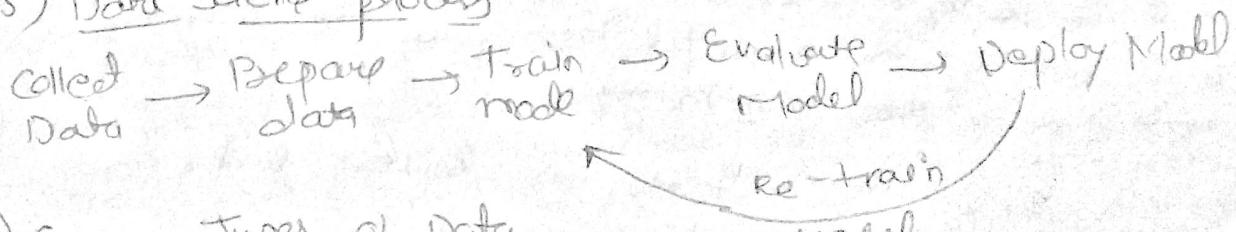
ML

→ A subcategory of AI that involves learning from data without being explicitly programmed

Deep learning

→ A subcategory of ML that uses a layered neural network architecture originally inspired by the human brain

5) Data Science process



6) Common Types of Data

Used in ML

- Numerical

- Time Series

- Categorical
- Text
- Image

→ gets all numerical in the end

7) Tabular data: that is data arranged in a data table, it consist row, column & cell

8) Scaling data: - it mean to transform it so that the value fit within some range or scale such as 0-100 or 0-10 or 0-1.

Eg: if we have an image - RGB whose range is 0-255 & we can scale down it into range of 0-1. The scaling process not affect algorithm output, it speed up the training process, because now algorithm handle has. less than or equal to 1

• Approach of Scaling

9) Standardization

$$\rightarrow \text{Mean} = 0 \text{ & Variance} = 1$$

10) Normalization

$$\rightarrow \text{range } [0, 1] \\ \rightarrow (x - x_{\min}) / (x_{\max} - x_{\min})$$

④ Encoding Categorical Data

1) Ordinal encoding

→ In ordinal encoding, we convert categorical data into

integer code ranging 0 to (no. of category - 1)

Eg:- Colour encoding $\frac{\text{Colour}}{0}$

$$\begin{array}{l} \text{Blue} \rightarrow 0 \\ \text{Green} \rightarrow 1 \\ \text{Red} \rightarrow 2 \end{array}$$

But we can't use it because it make confusion that red is more powerful than green the why we need 0 to red & 1 to green

→ one hot encoding: - If this we convert each other - cat value into a column.

| SKU | Red | Green | Blue |
|-------|-----|-------|------|
| 00121 | 0 | 0 | 0 |
| 00231 | 1 | 0 | 0 |
| 00441 | 0 | 0 | 1 |

- Algorithm understand by if two vectors are closer to each other, we can say text represent by two vector have similar meaning.

Eg:- quick fox lazy dog

| | | | | |
|-------------|------|------|------|------|
| [quick fox] | 0.32 | 0.23 | 0.0 | 0.0 |
| [lazy dog] | 0.0 | 0.0 | 0.12 | 0.23 |

- Whole pipeline to convert text
 - normalization
 - Feature extraction (Vectorization)
 - Algorithm

18) Two Perspective on ML

15) Computer science vs Statistical perspective

- Input feature & output feature
- Output = Program (Input Feature)
- A group of input variables is called Input vector
- We use program from input data to make prediction
- Eg:- we if take data in tabular form, row are entities
- columns are feature describe the property of entity

i) Output variable = $f(\text{Input variable})$
 Or
 Dependent = $f(\text{Independent variable})$
 Or
 $y = f(x)$

Eg:- price of a house is dependent variable that depend upon some variables like location, size

15) The tools for ML

o 1) Libraries :-

1) Scikit-learn 0) Tensorflow 0) PyTorch 0) Keras

o Development Environment -

- Jupyter notebook
- Azure notebook
- Azure DataBrick

- Visual Studio Code
- Visual Studio

o Cloud Service

Microsoft Azure ML

16) Libraries for ML

- Code Framework & Tools - python, pandas, Numpy, jupyter
- ML Deep learning - tensorflow, keras, pytorch

a) Visualization - Seaborn, matplotlib, bokeh, plotly

17) Cloud service in ML

Core asset management

- Dataset
- experiments/RUNS
- Pipeline
- Models
- Endpoint

Resource management

- Compute
- Environments
- Databases

- Dataset :- Define, version & monitor datasets used in ML
- Experiment/RUN :- Organize ML workload & keep track each task executed through the service.
- Pipeline :- Structural flows of tasks to make complex ML flow.
- Models :- Model artifact with support for versioning & deploy - map to production
- Endpt :- Expose real-time endpoints for scoring as well as pipeline for advanced automation.
- Compute :- Manage compute resource used by ML task.
- Environment :- Template for standardized environments used to create compute resources
- Database :- Data source connect to the service environment

18) Model vs Algorithm

- Model are the specific representation learned from data.
- Algorithm are the process of learning.
 $\text{Model} = \text{Algorithm(Data)}$

Q1) Linear Regression :-

- o Linear Regression is a algorithm that uses a straight line (or plane) to describe relationship b/w variables.

1) Simple linear Regression

$$y = mx + c \rightarrow \text{Intercept}$$

m slope

In ML, c is bias.

$$y = \beta_0 + \beta_1 x + \epsilon \quad (\text{simple linear regression})$$

2) Multiple Linear Regression

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n$$

Instead of line a plane is formed

- To train a linear regression model, we have to learn coeff & bias that best fit the data & minimize the error during prediction. To calculate error, we use cost function.
- Most commonly used cost function for linear regression is root mean square error (RMSE)

- i) The relation b/w input & output must be linear, if it is not linear we have to transform it.
- ii) To remove highly correlated input variable & reduce a correlation check among input variable.
- iii) Gaussian distribution.
- iv) Rescale data: we should first normalize & standardize the data
- v) Remove noise: Linear regression is very sensitive to noise & outliers in the data.

These are the things imp for preparation data

- formula for slope

$$m = \frac{\sum_{i=1}^n (x_i - \text{mean}(x)) \times (y_i - \text{mean}(y))}{\sum_{i=1}^n (x_i - \text{mean}(x))^2}$$

- To get intercept,

$$c = \text{mean}(y) - mx \text{mean}(x)$$

- root mean squared error

$$= \sqrt{\frac{1}{n} \sum_{i=1}^n (p_i - y_i)^2}$$

- Sklearn & Pytorch library there inner working for linear reg is already present.

2st learning function

$$y = f(x) + e \quad \begin{matrix} \rightarrow \text{irreducible error} \\ (\text{it is independent of } f(x) \text{ or input}) \end{matrix}$$

σ^2 is e due to don't have enough data or don't have enough data feature.

→ Model error measure how much the prediction made by the model is different from true output.

26. Parametric Vs Non-Parametric

Based on assumption about shape & structure of func they try to learn, ML algorithms can be divided into two category:

| 1) Parametric | 2) Non-parametric |
|--|--|
| i) Simplify the mapping to a known functional form (eg. linear regression) | • Not making assumption regarding the form of mapping b/w input & output (eg:- K Mean Neighbour (KNN)) |
| ii) Benefit :- Simple, Faster less Training Data | ii) Benefit :- High Flexibility, power, high performance |
| iii) Limitation:- Highly constrained, limited complexity, poor fit | iii) Limitation:- More Training Data, slower, overfitting training |

Dr. Classified ML vs Deep learning
learning algorithm and ML is not all for algorithm

• all steps during abortion due deep uterine endometrial changes resulting in normal uterine bleeding

Desp. leaving on board on David's Dredging

S. Green and S.

Deep learning - advantages
1) Suitable for high complexity

D) Difficult to explain
framed data

problem
2) Better accuracy, compared to class

2) require significant computational power

→ read NL
3) Better suggest for big data
4) Complex feature can be learned

卷之三

Classical ML advantage

- 1) More suitable for small datasets
- 2) Easier to implement and maintain
- 3) Can run on less powerful machines
- 4) Cheaper to perform
- 5) Does not require large compilation power

Difficult to learn large
changes
Require feature engineering
Difficult to learn
complex func.

Approach to Machine Learning

- 1) Supervised
 - 2) Unsupervised
 - 3) Reinforcement

Supervised learning

- Supervised learning — both the input and output variables are labeled. Common types are classification and regression.

- Output: e.g. labels
- Classification = output one categorical class
- Regression = Output new continuous numerical value
- Similarity reasoning = Learn from e.g. using a similarity measure how similar two objects are

- v) Feature learning - Learn to automatically discern the specific patterns or features from raw data.

v) Anomaly detection - A special form of classification, which learns from data labelled as normal / abnormal.

Unsupervised learning

- Unsupervised learning is my final hidden structure in mind, data

- 1) Learn from input data only
- 2) Clustering: Assign entities to clusters or groups
- 3) Feature learning: Features are learned from unlabeled data
- 4) Anomaly detection: Learn from unlabeled data, using assumption that majority of entities are normal.