



Dr. Vishwanath Karad
**MIT WORLD PEACE
UNIVERSITY** | PUNE
TECHNOLOGY, RESEARCH, SOCIAL INNOVATION & PARTNERSHIPS

Course Code	CET3009B			
Course Category	Professional Core			
Course Title	Data Science for Engineers			
Weekly Teaching Hrs. and Credits	L	T	Laboratory	Credits
	2	-	2	2 + 0 + 1
Pre-requisites: Linux Based Python Laboratory				
Course Objectives: <ol style="list-style-type: none">1. Knowledge: (i) To know fundamentals of data science and apply python concept for data analysis.2. Skills : (i) To learn basic concepts of statistics for data analysis. (ii) To learn data visualization tool and techniques for data analysis.3. Attitude: (i) To identify machine learning algorithm to solve real world problems.				
Course Outcomes: After completion of the course the students will be able to: - <ol style="list-style-type: none">1. Understand fundamentals of data science and python concepts for data analysis.2. Apply statistical concepts to solve real life problems.3. Apply appropriate machine learning algorithms to solve real world problems.4. Apply Visualization tool and techniques to find insights from real world data.				
Course Contents: <ol style="list-style-type: none">1. Introduction to Data Science2. Statistics for Data Science3. Machine Learning4. Data Visualization				
Laboratory Exercises / Practical: <ol style="list-style-type: none">1. Python Basic programming2. Data Preprocessing using Numpy and Pandas3. Data Preprocessing using Numpy and Pandas4. Basic Statistics using Python5. Simple Linear Regression6. Classification using Naive Bays7. Clustering Using K-Means8. Data Visualization using Python				

MIT-WPU
Approved by
Academic Council
23 JUL 2022
Date

Dr. Prasad Khandekar
Dean

Learning Resources:

Text Books:

1. Cathy O'Neil, Rachel Schutt, Doing Data Science, Straight Talk from The Frontline. O'Reilly, 2013
2. Applied Statistics and Probability for Engineers – By Douglas Montgomery
3. Jiawei Han, Micheline Kamber, Jian Pei, "Data Mining: Concepts and Techniques", 3rd Edition

Reference Books:

1. Foundations of Data Science by Avrim Blum, John Hopcroft, and Ravindran Kannan
2. Ward, Grinstein Keim, Interactive Data Visualization: Foundations, Techniques, and Applications. Natick: A K Peters, Ltd.
3. Glenn J. Myatt, Making sense of Data: A practical Guide to Exploratory Data Analysis and Data Mining, John Wiley Publishers, 2007.

Supplementary Reading:

https://swayam.gov.in/nd1_noc19_cs60/preview

Web Resources:

<https://nptel.ac.in/courses/106/106/106106179/>

Weblinks:

https://www.youtube.com/watch?v=MiiANxRHSv4https://www.youtube.com/watch?v=y8Etr3Tx6yE&list=PLYqSpQzTE6M_JcleDbrVyPnE0PixKs2JE&index=5

MOOCs:

<https://intellipaat.com/data-scientist-course-training/>

Pedagogy:

- PowerPoint Presentation
- Flipped Classroom Activity
- Project based Learning
- Jupyter notebook for coding

Assessment Scheme:

Class Continuous Assessment (CCA) (30 Marks)

Assignments	Mid Test	MCQ/Poster Presentation (Research Statement)/Active Learning
10	15	5

Laboratory Continuous Assessment (LCA) (30 Marks)

Understanding the Objectives	Understanding of Procedure and Initiatives	Experimental Skills	Oral
5	5	5	15

Term End Examination: Term end exam of 40 marks will be based on entire syllabus.

Dr. Prasad Khandekar
Dean




MIT-WPU
Approved by
Academic Council
23 JUL 2022
Date:

Theory Syllabus:

Unit	Contents	Workload in Hrs	
		Theory	Lab
1	Introduction to Data Science: Data Science Fundamentals: Types of Data, Data Quality, Data Science Life Cycle, Applications, Types of datasets, Python for Data Science: Pandas and Numpy, Matplotlib for data analysis, Data Pre-processing: Missing data handling, Data scaling and normalization, Feature extraction.	8	
2	Statistics for Data Science: Basic Statistics: Descriptive Statistics, Measures of Central Tendency: Mean, Median, Mode, Measures of Dispersion: Range, Variance, Standard Deviation, Measures of Position: Quartiles, Percentile, Z-score, Data transformation, Measure of Relationship: Covariance, Correlation, Basic Probability and Distribution, Hypothesis testing, Applying statistical concepts in Python.	9	
3	Machine Learning: Introduction to machine learning, Supervised and Unsupervised Learning, splitting datasets: Training and Testing, Regression: Simple Linear Regression, Classification: Naïve Bayes classifier and clustering: K-means, Evaluating model performance, Python libraries for machine learning.	9	
4	Data Visualization: Introduction to data visualization, challenges, Types of Data visualization: Bar charts, scatter plots, Histogram, Box Plots, Heatmap, Data Visualization using python: matplotlib, seaborn, Data Visualization tool: Tableau.	8	

Laboratory Assignments:

Assignment No	Title of the Assignment	Workload in Hrs.
1	<p>Attempt any 3</p> <ol style="list-style-type: none"> Write a python program to create a dictionary which contains student's names and marks. Iterate over the dictionary and apply below conditions to print their grades: <ol style="list-style-type: none"> Marks greater than or equal to 70 – Distinction Marks between 60-69 – First Class Marks between 50-59 – Second Class Marks between 40-49 – Pass Marks less than 40 - Fail Write a Python Program to create a 1D array of numbers from 0 to 9. 	2

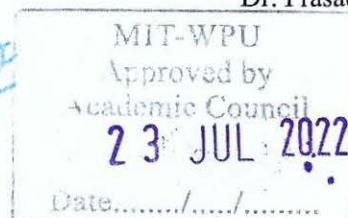
Dr. Prasad Khandekar

Dean

MIT-WPU
Approved by
Academic Council
23 JUL 2022
Date.....

	<p>iii. Write a NumPy program to create an array of all the even integers from 30 to 70.</p> <p>iv. Write a NumPy program to create a 3x4 matrix filled with values from 10 to 21.</p> <p>v. Write a NumPy program to compute the sum of all elements, sum of each column and sum of each row of a given array.</p>	
2	<p>Attempt any 3</p> <p>i. Write a python program to output a 3-by-3 array of random numbers following normal distribution Stack these arrays vertically: a = np.arange(10).reshape(2,-1) b = np.repeat(1, 10).reshape(2,-1)</p> <p>ii. Get the common items between two numpy arrays a = np.array([1,2,3,2,3,4,3,4,5,6]) b = np.array([7,2,10,2,7,4,9,4,9,8])</p> <p>iii. Create a series from a list, numpy array and dictionary Combine many series to make a data frame.</p> <p>iv. Create a normalized form of iris's sepal length whose values range exactly between 0 and 1 so that the minimum has value 0 and maximum has value 1. Input: url = 'https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data' sepal length = np.genfromtxt(url, delimiter=',', dtype='float', usecols=[0]) Hint: Apply Min-Max Scalar formula</p>	2
3	<p>Load Data and perform Data Pre-processing. Input: df = pd.read_csv (https://raw.githubusercontent.com/selva86/datasets/master/Cars93_miss.csv)</p> <p>i. Read a csv file to create a data frame and print top records.</p> <p>ii. Check if there are any missing values in the data.</p> <p>iii. Drop null values / Impute the missing values with mean / median.</p> <p>iv. Import 'crim' and 'medv' columns of the BostonHousing dataset as a dataframe and get the n rows, n columns, datatype, summary stats of each column of a dataframe.</p> <p>v. Which manufacturer, model and type has the highest Price?</p> <p>vi. How to create one-hot encodings of a categorical variable.</p>	2
4	<p>Understanding Statistical concepts in Python. (Attempt any 3)</p> <p>i. The average test scores are given: test scores: 83,85,87,89,91,93,95,97,99,100. Find Mean, Median,</p>	2

Dr. Prasad Khandekar
Dean



	<p>Variance, Standard deviation of the data. Show the information on the bell curve.</p> <p>ii. Consider given product price data: price_data=[13,43,54,34,40,56,34,61,34,23]. Find Range, 25th Percentile and IQR.</p> <p>iii. A person tries to analyse the last 12 months interest rate of the investment firm to understand the risk factor for the future investment. The interest rates are: 12.05%, 13%, 11%, 18%, 10%, 11.5%, 15.08%, 21%, 6%, 8%, 13.2%, 7.5%.</p> <table><thead><tr><th>Months (One Year)</th><th>Interest Rate (%)</th></tr></thead><tbody><tr><td>April</td><td>12.05</td></tr><tr><td>May</td><td>13</td></tr><tr><td>June</td><td>11</td></tr><tr><td>July</td><td>18</td></tr><tr><td>August</td><td>10</td></tr><tr><td>September</td><td>11.5</td></tr><tr><td>October</td><td>15.08</td></tr><tr><td>November</td><td>21</td></tr><tr><td>December</td><td>6</td></tr><tr><td>January</td><td>8</td></tr><tr><td>February</td><td>13.2</td></tr><tr><td>March</td><td>7.5</td></tr></tbody></table> <p>iv. Calculate Skewness and Kurtosis and comment on it.</p> <p>v. Hypothesis Testing</p> <p>a. Consider below data and tests whether a data sample has a Gaussian distribution by formulating hypothesis test</p> <p>b. data = [0.873, 2.817, 0.121, -0.945, -0.055, -1.436, 0.360, -1.478, -1.637, -1.869]</p>	Months (One Year)	Interest Rate (%)	April	12.05	May	13	June	11	July	18	August	10	September	11.5	October	15.08	November	21	December	6	January	8	February	13.2	March	7.5	
Months (One Year)	Interest Rate (%)																											
April	12.05																											
May	13																											
June	11																											
July	18																											
August	10																											
September	11.5																											
October	15.08																											
November	21																											
December	6																											
January	8																											
February	13.2																											
March	7.5																											
5	<p>Write a python program to predict the height of a person providing his age using the trained model to the highest achievable accuracy using available data.</p> <p>Perform following steps:</p> <p>i. Importing the dataset. Link of Data.</p> <p>ii. Perform exploratory analysis of the data: Print features, Shape, Size, labels, head records, data types, outliers etc.</p> <p>iii. Data Cleaning.</p> <p>iv. Build the Model and Train it.</p> <p>v. Make Predictions on Unseen Data.</p>	2																										

Dr. Prasad Khandekar
Dean





TECHNOLOGY, RESEARCH, SOCIAL INNOVATION & PARTNERSHIPS

vi. Analyze the performance of the model.

6

Write a python program to build a model to classify the type of cancer. The data has two types of cancer classes: malignant (harmful) and benign (not harmful). Perform following steps:

- Load the Data (The dataset is available in the scikit-learn library).
- Exploring Data: Prints features, Shape, Size, labels, head records, data types, outliers etc.
- Split the data into train and test set.
- Select the classification model.
- Fit the model on train Data.
- Predict the outcome on test data.
- Evaluate the performance of the model: Confusion matrix, accuracy, F1, Precision, Recall.
- Check of Tuning Hyperparameters of the model to improve performance.

4

7

Write a python program to perform Clustering: We have the data for workout as below.

date	distance_km	duration_min	delta_last_workout	day_category
10/17/17	4.3	21.58	1	0
11/04/17	1.9	9.25	18	1
11/18/17	1.9	9.0	14	1
11/23/17	1.9	8.93	5	0
11/28/17	2.3	11.94	5	0
11/29/17	2.8	14.05	1	0

To keep track of your performance you need to identify similar workout sessions. Clustering can help you group the data into distinct groups, guaranteeing that the data points in each group are similar to each other. Perform following steps:

- Load the Data
- Data Exploratory Analysis: Pair Plot and Distance versus workout duration, distance versus duration with the number of days, correlation (Scatter plot) to get idea about correlation between different features.
- Select K-means clustering for model and get the clusters.
- Evaluate the performance of the model.

4




MIT-WPU Dr. Prasad Khandekar
Approved by
Academic Council
23 JUL 2022
Date...../...../.....



Dean

8	<p>Download company sales data and perform following operations (Attempt any 5).</p> <ol style="list-style-type: none"> Read Total profit of all months and show it using a line plot. Generate above plot with following style properties <ol style="list-style-type: none"> Line Style dotted and Line-color should be red Show legend at the lower right location. X label name = Month Number Y label name = Sold units number Add a circle marker. Line marker color as read Line width should be 3 Read the total profit of each month and show it using the histogram to see the most common profit ranges. Calculate total sale data for last year for each product and show it using a Pie chart. Read all product sales data and show it using the stack plot. Read all product sales data and show it using a multiline plot. Display the number of units sold per month for each product using multiline plots. (i.e., Separate Plotline for each product). Read toothpaste sales data of each month and show it using a scatter plot. Read face cream and facewash product sales data and show it using the bar chart. 	4
---	-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	---





Dr. Prasad Khandekar
Dean

