# Lahore University of Management Sciences

## CS 334/EE 402 – Principles and Techniques of Data Science
Spring 2024
<mark>Subject to Change</mark>

| Instructor | Dr. Mobin Javed | | | |
|---|---|---|---|---|
| Room No. | SBASSE 9-G10A | | | |
| Email | mobin.javed@lums.edu.pk | | | |
| Telephone | 3338 | | | |
| TAs | TBA | | | |
| TA Office Hours | TBA | | | |
| Course URL (if any) | lms.lums.edu.pk | | | |
| **Course Basics** | | | | |
| Credit Hours | 3 | | | |
| Lecture(s) | Nbr of Lec(s) Per Week | 2 | Duration | 75 mins |
| Recitation/Lab (per week) | Nbr of Lec(s) Per Week | -- | Duration | |
| Tutorial (per week) | Nbr of Lec(s) Per Week | -- | Duration | |
| **Course Distribution** | | | | |
| Core | No | | | |
| Elective | Yes | | | |
| Open for Student Category | All | | | |
| Close for Student Category | None | | | |

| COURSE DESCRIPTION |
|---|

This is an introductory-level Python-based course in data science to prepare students for scientific work, as well as advanced courses in data mining and machine learning. It is a hands-on course and involves data analysis work.

The first half will focus on the fundamentals. We will start with descriptive statistics, develop an understanding of the biases in working with data, and learn and practice exploratory data analysis. The second half will focus on drawing inferences from data -- we will talk about setting up controlled experiments, hypothesis testing, and the foundational concepts of statistical and machine learning. We will also spend a major part of the course learning about data engineering, i.e., tools and techniques for data collection, data storage and querying, and working with big data.

| COURSE PREREQUISITE(S) | |
|---|---|
| | - CS100: Computational Problem Solving |
| | - MATH120: Linear Algebra with Differential Equations |
| | (Exceptions possible with prior instructor approval) |
| | - Basic knowledge of Probability and Statistics is assumed |
| | - Basic knowledge of Python programming language is assumed |

# Lahore University of Management Sciences

| COURSE OBJECTIVES | |
|---|---|
| | The goal of this course is to train students to become good data scientists. The students will build a foundation in drawing inferences from data, which serves as a preparation for advanced courses in data mining and machine learning. |

| Learning Outcomes | |
|---|---|
| | • Learn how to conduct sound data analysis<br>• Learn how to describe a given dataset and assess its quality<br>• Understand issues in experiments involving active data collection and develop the discipline of maintaining meta-data<br>• Learn how to build data pipelines (collection, cleaning, EDA, modeling, evaluation, results) for "repeatable" work<br>• Become well-versed with tools and technologies for data analysis (e.g., Pandas, Spark, scitkit-learn, R)<br>• Learn the theory behind drawing inferences from data<br>• Learn how to communicate results effectively |

## Grading Breakup and Policy

HWs:                    30%
Quizzes:             25%
Project:               15%
Final:                   25%
Labs:                    0%
Class Participation:  5%

**HWs:**
We will have four homeworks during the course of the semester, one corresponding to each of the first four modules.
**Optional Labs:**
Labs will be released to help you practice the concepts in the lecture. These labs are optional and will *not* count towards your grade.
**Quizzes:**
We will have a total of eight announced quizzes. Out of these best six will count towards your final grade (N-2 policy).
The quizzes will be held on LMS during lecture timings. The exact date of each quiz will be announced in lecture and on Piazza.
**Projects:**
Projects must be done in teams of 3-4.
The spirit behind the project is to give you some experience conducting a data science project end-to-end. This includes thinking about what question(s) to answer, gathering the right dataset, data quality assessment, EDA, drawing inferences from data, and building models. In addition, we want you to get some practice in storytelling and communicating your data science work clearly and crisply. Project is open ended by design, you are expected to come up with ideas on what to do and how to do. We have designed four check-points during the semester to help guide your exploration. More information on these will be released during the course of the semester.

## Examination Detail

| | |
|---|---|
| Midterm Exam | Yes/No:  No<br>Combine Separate:<br>Duration:<br>Preferred Date:<br>Exam Specifications: |

# Lahore University of Management Sciences

| Final Exam | Yes/No: Yes<br>Combine Separate:<br>Duration: Exam<br>Specifications: |
|---|---|

# Lahore University of Management Sciences

| Lec # | Topics | Assessments |
|-------|--------|-------------|
| 1 | Overview of Data Science: Untangling the Data Science Process<br>• What is Data Science? Why Data Science? Why Now?<br>• Data Science Lifecycle<br>• Overview of course modules | **Lab-0 Release (Setup)** |
| **Module 1: Descriptive Statistics, Data Acquisition, and Tools** | | **Quiz Dates TBA** |
| 2 | Descriptive Statistics, Deceptive Descriptions, Important Distributions, Choosing Unit of Analysis, Data Acquisition, Sampling, and Sources of Bias | |
| 3 | Data Manipulation Using Pandas – I (Lecture + Lab)<br>• Jupyter Notebooks<br>• Introduction to Pandas | **Lab-1 Release (Optional)** |
| 4 | Data Manipulation Using Pandas – II (Lecture + Lab)<br>• Data Aggregation<br>• Case Study | **HW-1 Release** |
| **Module 2: Data Cleaning, Exploratory Data Analysis, and Visualization** | | **Quiz Dates TBA** |
| 5 | Data Cleaning and Exploratory Data Analysis (EDA) – I<br>• Transforming Data for Ease of Analysis<br>• Common Data Anomalies | |
| 6 | EDA – II and Data Visualization and Transformations – I<br>• Structure, Granularity, Scope, Temporality, Faithfulness of Data<br>• Visual Representations of Data: A Way of Amplifying Cognition | |
| 7 | Data Visualization and Transformations – II<br>• Principles of Sound Visualizations<br>• Smoothing and Transformations | |
| 8 | Analyzing Text Data<br>• Text Mining and String Manipulation<br>• Regular Expressions | |
| 9 | Databases and SQL<br>• Relational Databases<br>• SQL Queries | **HW-2 Release** |
| **Module 3: Experiments, Causality, and Foundations of Statistical Inference** | | **Quiz Dates TBA** |
| 10 | Experiments, Observational Studies, and Causal Inference – I<br>• Association, Causality, and Data<br>• Randomized Control Trials and Case Study | |
| 11 | Causal Inference – II<br>• Estimating the Counterfactual<br>• Rubin's Causal Model of Potential Outcomes | |
| 12 | Causal Inference – III<br>• Causal Inference from Observational Studies (RDD, IV etc)<br>• Causal DAGs | |

# Lahore University of Management Sciences

| 13 | Statistical Inference and Hypothesis Testing<br>  • Sampling, Assessing Models, and Comparing Distributions<br>  • Hypothesis Testing and P-values<br>  • Case Study | |
|---|---|---|
| 14 | A/B Testing and Permutation Tests<br>  • Comparing Samples and Case Study | **HW-3 Release** |
| 15 | Bootstrap Sampling, Central Limit Theorem, Confidence Intervals<br>  • Repeated Random Sampling<br>  • Samples Averages<br>  • Variability and Bounds | |
| **Module 4: Machine Learning** | | **Quiz Dates TBA** |
| 16 | Models & Estimation<br>  • What is a Model?<br>  • Modeling Process and Loss Functions | |
| 17 | Optimization and Gradient Descent<br>  • Convexity and Gradient Descent<br>  • Stochastic Gradient Descent | |
| 18 | Regression and Linear Models<br>  • Simple Linear Regression<br>  • Ordinary Least Squares | |
| 19 | **Pakistan Day Holiday (Mar 23)**<br>Class Rescheduled to Apr 25: Industry Guest Lecture | |
| 20 | Feature Engineering<br>  • What is Feature Engineering?<br>  • Pitfalls | |
| 21 | Fundamental Challenges in Learning: Bias-Variance Tradeoff<br>  • Risk and Cost Minimization<br>  • Models Bias and Variance<br>  • Bias-Variance Tradeoff | **HW-4 Release** |
| 22 | Regularization and Cross-Validation<br>  • Train-Test Split | |
| 23 | Classification and Logistic Regression<br>  • Logistic Function<br>  • Stochastic Gradient Descent | |
| **Module 5: Big Data and Ethics** | | **Quiz Dates TBA** |
| 24 | Big Data Processing - I<br>  • Storage: Distributed File Systems<br>  • MapReduce and Spark | **Lab-2 Release (Optional)** |
| 25 | Big Data Processing – II & Ethics in Data Science | |

**Lahore University of Management Sciences**

| 26 | Guest Lecture: Ethics in Data Science | |

| Textbook(s)/Supplementary Readings |
| --- |
| There is no specific textbook for this course. Following are recommended readings:<br><br>(i) Principles and Techniques of Data Science<br>　　The textbook for DS100 at UC Berkeley<br><br>(ii) Doing Data Science<br><br>(iii) Naked Statistics: Stripping the Dread from the Data<br>　　A short good read on the fundamentals of Probability and Stats |